# Common Pitfalls in Computer Vision & AI Projects (and How to Avoid Them)

*A pragmatic playbook with lessons learned from real-world examples.*

Satya Mallick, Ph.D. | CEO

**B:GVISION**

# Table Of Contents

# Introduction

[Over 85% of AI initiatives fail to deliver their intended business value](#).

Computer Vision and AI projects hold immense promise; however, industry statistics reveal a sobering reality.

This playbook emerges from the trenches of actual project implementations—a compilation of hard-won lessons from successful deployments, costly failures, and everything in between.

Unlike theoretical guides, this document focuses exclusively on practical pitfalls that repeatedly surface across organizations of all sizes, from ambitious startups to Fortune 500 enterprises.

## Who Should Use This Guide

This playbook serves technical leaders, project managers, data scientists, engineers, and executives who are planning, executing, or evaluating Computer Vision and AI initiatives.

## Common Pitfalls & How to Avoid Them

Pitfalls lurk in every phase of an AI product cycle—from whiteboard to production logs. Let's name them, shame them, and eliminate them.

### 1. Fuzzy problem definition and success criteria

**Pitfall.** Teams jump into modeling without a crisp, testable problem statement.

**Avoid it.**

- Write a one-page PRD: who is the user, what decision the model informs, and what happens if it's wrong.

- Translate business goals → measurable success metrics (e.g., "cut manual review time by 40% at ≤1% false-positive rate").

- Pre-commit to a **single primary metric** and a **few key guardrail metrics** (latency, cost, fairness, safety, and privacy).

**Checklist.** Problem statement, target users, scope boundaries, primary metric & target, guardrails, acceptance test, rollback plan.

## 2. Unrepresentative Training data

**Pitfall.** Training data doesn't represent production reality—missing edge cases, long-tail classes, or operational conditions (lighting, motion blur, sensor differences).

**Avoid it.**

- Inventory real-world conditions. Turn them into **data requirements** (percentages per condition, per geography, per device).

- Collect deliberately: stratify by conditions; use **hard-negative mining**; synthesize sparingly to fill gaps; probe for **open-set** cases.

- Build an "edge-case suite" that is versioned and always evaluated.

**Case study.** Vision systems routinely fail under distribution shift—for example, face recognition accuracy dropped widely when masks appeared in 2020; production reality changed faster than training sets (see sources at the end).

## 3. Label quality, leakage, and splits

**Pitfall.** Noisy labels, inconsistent guidelines, and leakage (e.g., adjacent video frames or the same subject in train and test) inflate offline metrics.

**Avoid it.**

- Write **annotation guidelines** with concrete examples; run pilot labeling with adjudication.

- Track **inter-annotator agreement**; spot-check with a gold set.

- Split **by identity/scene/session** to prevent leakage; for medical or industrial data, split by **patient/part/lot/site**.

- Keep a small, untouched **holdout set** for go/no-go.

**Quick test.** If a simple random-forest model on metadata performs suspiciously well, you likely have leakage.

# 4. Focusing on a model-centric approach instead of a data-centric approach

**Pitfall.** Chasing architectures, tricks, and hyper-params while neglecting data coverage, label policy, and sampling. This yields brittle wins (if any) that don't transfer to holdout or production, and slows time-to-value.

**Avoid it.**

- Treat **data as the product**: write a data PRD with coverage targets by condition/device/site and per-class quotas.
- Run **data error analysis** every sprint: slice metrics, review top failures, tag error modes (e.g., glare, motion blur, near-duplicates, boundary cases), and file **data bugs**.
- Fix the taxonomy & guidelines; run **re-label sprints** with adjudication; **hard-negative mine** and deduplicate before splitting.
- Use a **small, fast model** (or frozen encoder + linear probe) to iterate data quickly; only scale model capacity after data metrics plateau.
- Track **data KPIs**: label-noise rate, IAA, duplicate rate, edge-case coverage, OOD/unknown rejection rate.

**Signals you're stuck in model-centric mode.**

- Frequent architecture swaps yield <1–2pp on holdout; gains vanish in shadow/canary.
- Big improvements occur when changing dataset source or split—not when changing the model.
- PRs show lots of code tweaks but few labeled data patches; failure modes repeat release to release.

**Case Study**

| Approach | Steel Defect Detection | Solar Panel |
|---|---|---|
| Baseline | 76.2% | 75.68% |
| Model-centric | 0% | 0.04% |

| Approach | Steel Defect Detection | Solar Panel |
|---|---|---|
| Data-centric | 16.9% | 3.06% |

*Table from slides shared by Dr. Andrew Ng.*
*https://www.youtube.com/live/06-AZXmwHjo?feature=shared&t=386*

# 5. Shortcut learning & spurious correlations

**Pitfall.** Models latch onto proxies (logos, backgrounds, timestamp overlays) instead of the object of interest.

**Avoid it.**

- Use **causal probes**: occlusion tests, counterfactual edits, and Grad-CAM style saliency sanity checks.

- Curate **counter-examples** (same label, different context; different label, same context).

- During collection, randomize backgrounds, cameras, and operators.

# 6. Misaligned metrics and thresholds

**Pitfall.** Optimizing AP/accuracy while the business runs on the cost of errors, per-class SLAs, or human-review throughput.

**Avoid it.**

- Convert business costs to a **utility curve**; choose operating points by cost, not just F1.

- Report **per-segment** metrics (device, site, skin tone, lighting, class frequency).

- Track **calibration** (ECE/Brier); use Platt/temperature scaling if needed.

**Template.** For each release, publish PR/ROC curves, confusion matrices, calibration plots, and recommended thresholds per segment.

# 7. Underpowered baselines & over-engineering

**Pitfall.** Starting with a giant foundation model without a simple baseline.

**Avoid it.**

- Establish **strong baselines** (resized classical features, a small CNN, and a linear probe on frozen encoders).

- Only add complexity if it wins on **primary + guardrail** metrics and on the **holdout** set.

# 8. Reproducibility & experiment drift

**Pitfall.** "Can't reproduce last week's win." Data changed, seeds changed, or augmentation/code drifted.

**Avoid it.**

- Version data, code, configs, and model artifacts (e.g., DVC + Git + Weights & Biases/MLflow).

- Fix seeds; log all hyperparams and environment; store training/val indices.

- Use **ablations** and **checklists** for each experiment.

# 9. Ignoring deployment constraints

**Pitfall.** A model that's great on a workstation but unusable on edge devices or in a high-QPS service.

**Avoid it.**

- Define **SLOs** up front: latency (p50/p95), memory, power, model size, cold-start, batch sizes.

- Plan the optimization path: **quantization**, **pruning**, **distillation**, operator fusion, and hardware-aware architectures.

- Profile early on target hardware.

## 10.  Monitoring, drift, and silent failures

**Pitfall.** "It shipped and then got worse." No telemetry, no drift detection, no human-in-the-loop.

**Avoid it.**

- Log inputs, scores, decisions (with privacy controls). Track **data & concept drift** and **calibration drift**.

- Run **shadow** or **canary** deployments; set auto-rollback triggers.

- Maintain a **review UI** and feedback loop; sample and relabel difficult cases via **active learning**.

## 11.  Fairness, safety, and societal impact

**Pitfall.** Uneven performance across subpopulations (skin tone, age, gender), and harmful failure modes in safety-critical contexts.

**Avoid it.**

- Define fairness guardrails and report segmented metrics.

- Stress-test under new conditions (masks, PPE, occlusions).

- Use human oversight where the risk of harm is high; document with **model cards** and **data sheets**.

**Case studies.**

- Public tests showed [high-profile face recognition mismatches in 2018](); later evaluations also found big [accuracy drops for masked faces in 2020–2022]().

- [Dermatology AI systems trained on light-skin images underperform on darker skin](); targeted fine-tuning narrows the gap.

## 12.  Security & adversarial robustness

**Pitfall.** Overlooking adversarial or opportunistic attacks ([sticker attacks on traffic signs](), adversarial patches, data-poisoning, prompt injection for VLMs).

**Avoid it.**

- Threat-model the system; add **input sanitation**, **rate limits**, and **ensemble checks**.

- Use **confidence thresholds** and **abstention**; detect OOD.

- Test with **physical-world adversarial suites** and randomized augmentations; monitor for distribution anomalies.

## 13. Human factors & UX

**Pitfall.** Automation bias (over-trust) or alert fatigue (under-trust). Poor handoff to human reviewers.

**Avoid it.**

- Design for **"fail-gracefully"**: show uncertainty, provide explanations suitable for the operator, and allow easy override.

- Measure **end-to-end** time-to-decision, not just model latency.

## 14. Privacy, IP, and compliance

**Pitfall.** Using data without proper rights or violating biometrics/privacy laws.

**Avoid it.**

- Maintain **data provenance** and licenses; document consent.

- For biometrics, review BIPA/GDPR/CCPA obligations; minimize retention; hash/pseudonymize IDs; enable deletion.

- Red-team prompts and outputs for sensitive content if using VLMs.

## 15. Integration & change management

**Pitfall.** The model works, but the workflow doesn't—no APIs, no SLAs, no process changes.

**Avoid it.**

- Own the **end-to-end** pipeline: data → model → service → review UI → analytics.

- Provide clear APIs, SLAs, and **runbooks**; train end users; measure adoption.

## 16. ROI illusions

**Pitfall.** Counting model accuracy as ROI. Underestimating ops/labeling/infra costs and the cost of human review.

**Avoid it.**

- Build a **unit-economics model** (labeling $/image, inference $/k images, human review minutes, expected error cost).

- Stage-gate funding by business impact; pilot with a **shadow test** and compare to manual baseline.

# "Pre-Flight" Checklist
## *(print and tape to your monitor)*

1. One-page PRD with primary metric & target.

2. Data specification with coverage targets by condition/segment.

3. Annotation guidelines + gold set + agreement tracked.

4. Data-centeric > Model-centric

5. Leakage-proof splits (by ID/scene/site).

6. Baselines established; ablation plan.

7. Edge-case test suite (versioned).

8. Guardrails: fairness segments, safety terms, privacy constraints.

9. Calibration check; operating thresholds by segment.

10. Deployment SLOs (latency, memory, cost) and profiling on target.

11. Model card + data sheet + rollback/runbook.

12. Monitoring plan (drift, quality, alerts).

13. Human-in-the-loop workflow defined.

14. Security review & adversarial test plan.

15. ROI model and stage-gates.

16. Shadow/canary plan and success criteria.

# Templates you can copy

**Model Card (skeleton)**

- **Intended use:** ...

- **Not intended for:** ...

- **Primary metric & target:** ...

- **Segments & guardrails:** ... (e.g., device/site/skin-tone)

- **Training data summary:** ...

- **Evaluation protocol:** splits, leakage checks, edge-case suite

- **Calibration & thresholding:** ...

- **Safety/fairness review:** ...

- **Deployment constraints:** latency/memory/power

- **Monitoring & retraining triggers:** ...

**Evaluation Plan (one-pager)**

- **Holdout definition:** ...

- **Edge-case suite:** ...

- **Metrics:** primary, guardrails, cost curve

- **Release gate:** numbers + qualitative checks (saliency, counterfactuals)

**Annotation SOP**

- Guidelines + examples → Pilot → IAA target → Adjudication → Ongoing QA → Drift re-training loop

## Final advice

- Treat **data as code** and **deployment as design**, not an afterthought.

- Prefer **abstention over confident mistakes** in high-risk contexts.

- If you can't explain how the model improves the business metric **and** how you'll know when it degrades, you're not ready to ship.

## Sources for the case studies

- Big Vision AI Consulting Case Studies

- ACLU's 2018 test of Amazon Rekognition (28 false matches of U.S. lawmakers). American Civil Liberties Union

- NIST FRVT reports on face recognition accuracy with masks (2020–2022 updates). NIST PagesNIST Publications

- Dermatology AI performance gaps on darker skin tones (Daneshjou et al., 2022; Groh et al., 2024). PubMedNature

- Physical-world adversarial stickers that fool traffic-sign classifiers (Eykholt/Evtimov et al., 2017). arXiv

- NTSB report on the 2018 Uber ATG fatal crash in Tempe (perception/ADS system-level failures). NTSB

- Google Photos' 2015 "gorillas" labeling incident and follow-on reporting. WIRED, The Guardian, eKathimerini

- Data-centric > Model-centric

## Contact

You can contact us by sending us an email at contact@bigvision.ai.