# AI Testing Playbook

From understanding the AI App to confidently test it

Dana Andrei – Testing Delivery Manager

# Summary

1. Why AI needs a new testing approach

2. Basic AI Concepts
   - AI Modalities
   - AI Systems
   - AI Capability model
   - AI Risks Considerations

3. AI Testing Playbook

# AI Testing Playbook

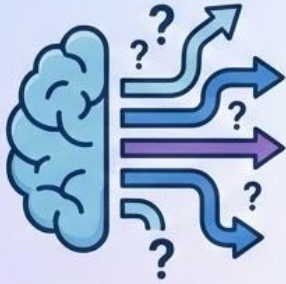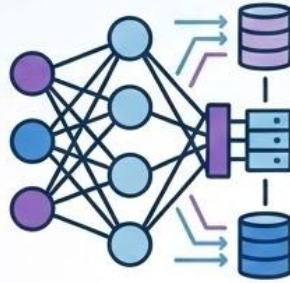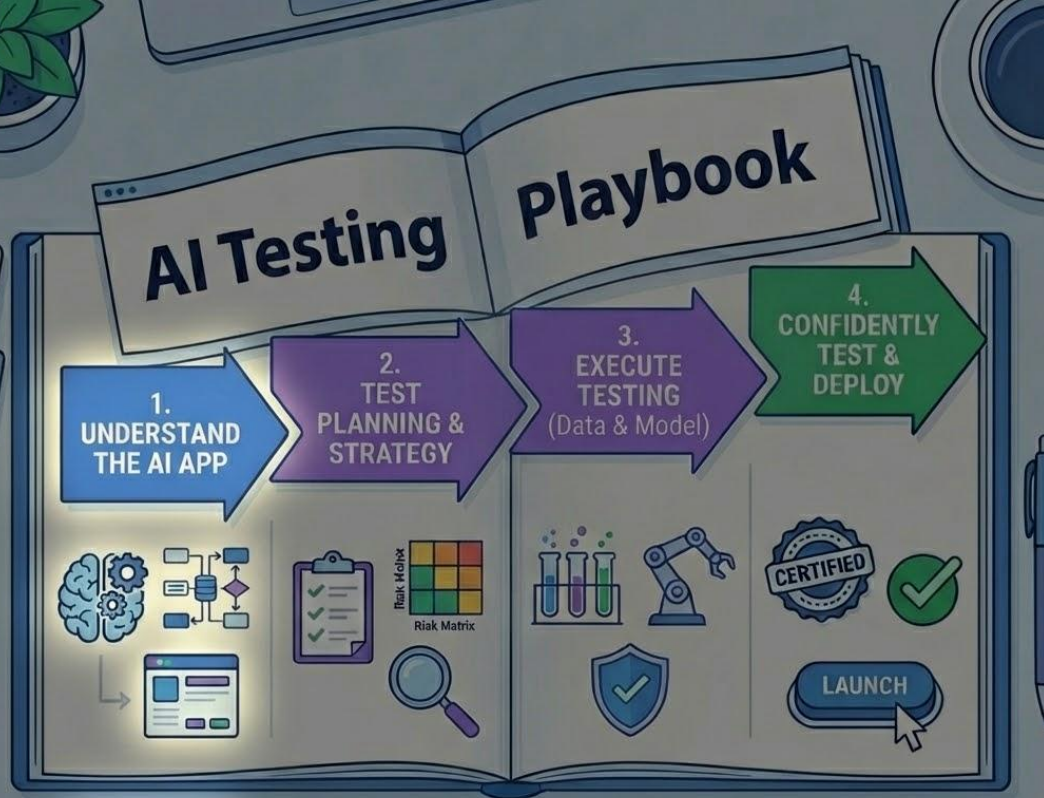| 1. UNDERSTAND THE AI APP | 2. TEST PLANNING & STRATEGY | 3. EXECUTE TESTING | 4. CONFIDENTLY TESTED & DEPLOY |

# Why AI needs a new testing approach?



AI is probabilistic, not deterministic; the same input can produce different outputs.



Traditional testing assumes fixed rules and repeatable outcomes. AI (especially LLMs and ML models) behaves based on patterns, patterns, probabilities, and training data.



Quality

Risk

Reasoning

This makes pass/fail testing insufficient — you must evaluate quality, risk, and reasoning rather than strict correctness, not deterministic — the same input can produce different outputs.

# AI Testing Playbook

# Basic AI Concepts

**AI Modalities**

**AI Systems**

**AI Capability Model**

**AI Risk Considerations**

# AI Modality

**Definition:** An AI modality refers to the type of input data or signal that an artificial intelligence system processes, analyzes, or generates. It defines the sensory or informational channel the AI "understands" and acts upon. Each modality requires specific models, architectures, and testing approaches.

## Key Points

A modality represents how the AI perceives the world: text, speech, image, video, sensor, tabular data, etc.

Some AI systems are single-modality (e.g., NLP-only chatbots), while others are multimodal, combining several modalities (e.g., GPT-4-Vision combining text + images).

Choosing the right modality is crucial for system design, training data, evaluation, and risk management.

# Practice Quiz: Match the AI Modality to the Risk

| AI Modality | Risk |
|---|---|
| 1 Text | A reinforcing stereotypes in job recommendations |
| 2 Image | B responding in wrong tone |
| 3 Structured Data | C misclassifying images |
| 4 Voice | D misinterpreting speech |

# Example of AI Modality

| AI Modality | Description | Real-World Examples | Primary Risk (What Can Go Wrong) |
|---|---|---|---|
| Text / NLP | Processes and generates natural language (written text) | ChatGPT, summarization tools, sentiment analysis, legal document analysis | **Hallucinations & misinformation** – generates fluent but incorrect or fabricated content |
| Speech / Audio | Processes spoken language, sound events, and audio signals | Siri, Alexa, call center transcription, voice biometrics | **Misrecognition & accent bias** – fails for certain accents, noise conditions, or speakers |
| Image / Vision | Processes visual inputs like images, photos, and diagrams | Face ID, Google Lens, medical imaging AI, factory defect detection | **False positives / negatives** – misidentification leading to incorrect decisions |
| Video | Processes sequential visual frames with or without audio | Surveillance analytics, TikTok recommendations, sports video analysis | **Context misinterpretation over time** – incorrect conclusions from partial or evolving scenes |
| Multimodal (Text + Image / Video / Audio) | Combines multiple data types for richer understanding | GPT-4 Vision, Document AI (OCR + text), DALL·E with image prompts | **Cross-modal hallucination** – invents relationships between modalities that don't exist |
| Tabular / Structured Data | Works on rows/columns of structured datasets | Credit scoring, predictive maintenance, sales forecasting | **Hidden bias & unfair decisions** – biased outcomes masked by "objective" numbers |
| Time-Series / Sequential Data | Works with sequential or temporal data | Stock prediction, sensor monitoring, IoT anomaly detection | **Concept drift** – model degrades silently as patterns change over time |
| Graph / Network Data | Works with nodes and edges representing relationships | Fraud detection, social network analysis, recommendation graphs | **Amplification of existing biases** – reinforces unfair or harmful network structures |
| 3D / Spatial / Point Cloud | Works with 3D structures, LIDAR, or spatial mapping | Self-driving car LIDAR, AR/VR spatial apps | **Spatial misalignment** – incorrect depth or distance interpretation causing unsafe actions |
| Reinforcement / Interaction Data | Learns through trial-and-error via environment feedback | AlphaGo, robotic manipulation, game AI | **Reward hacking** – optimizes for the wrong goal in unintended ways |
| Sensor / IoT Data | Reads physical sensor streams (temperature, motion, etc.) | Smart thermostats, industrial monitoring, wearables | **Sensor noise & failure propagation** – bad signals lead to bad decisions at scale |
| Genomic / Bioinformatics Data | Works with genetic or molecular sequences | Cancer mutation prediction, CRISPR design AI | **Overconfidence in probabilistic predictions** – high impact medical errors |
| Hybrid / Neuro-Symbolic Inputs | Combines symbolic rules with ML inputs | Medical decision support with guidelines + patient data | **Logic–model inconsistency** – rules and learned behavior contradict each other |

# Practice Quiz: Match the AI System Type to its Definition

| AI System Type | | Definition | |
|---|---|---|---|
| 1 | generative AI | A | performs goal-driven actions using memory or planning |
| 2 | autonomus AI | B | creates new content like text, summaries, audio, images, or code |
| 3 | agentic AI | C | acts independently in real environments |
| 4 | predictive AI | D | deterministic logic, no learning |
| 5 | reactive AI | E | predicts future outcomes or behaviours from historical/user data |

# AI Capability Maturity Model

Increasing Maturity, Capability, and Risk Profile

**Level 7: Autonomous Systems, Safety-critical AI**
Acts independently in real environments.
e.g., Drones, ADAS systems

**Level 6: Agentic AI**
Plans, executes, and adapts actions.
e.g., Auto-resolving support agents

**Level 5: Multimodal AI**
Reasons across multiple input types.
e.g., Image+text analysis

**Level 4: Generative AI (LLMs, image models)**
Produces novel content.
e.g., ChatGPT, image generation models

**Level 3: Context-Aware AI, Memory-limited AI, NLP systems**
Uses recent history or context.
e.g., Chatbots with session memory

**Level 2: Predictive Models, Recommenders**
Learns patterns from historical data.
e.g., Credit scoring, demand forecasting

**Level 1: Reactive, Rule-Based Automation**
Deterministic logic, no learning.
e.g., Spam rules, IVR menus

# AI Capability Maturity Model

| Maturity Level | Systems | Capability Description | Examples | Primary Testing Focus | Risk Profile |
|---|---|---|---|---|---|
| Level 1 | Reactive, Rule-Based Automation | Deterministic logic, no learning | Spam rules, decision trees, IVR menus | Functional correctness, edge cases, rule conflicts | Low |
| Level 2 | Predictive Models, Recommenders | Learns patterns from historical data | Credit scoring, churn prediction, demand forecasting | Data quality, bias, accuracy drift | Medium |
| Level 3 | Context-Aware AI, Memory-limited AI, NLP systems | Uses recent history or context | Chatbots with session memory, voice assistants | Context retention, hallucination, partial failures | Medium–High |
| Level 4 | Generative AI (LLMs, image models) | Produces novel content | ChatGPT, Copilot, summarization AI | Output quality, safety, non-determinism, evaluation rubrics | High |
| Level 5 | Multimodal AI | Reasons across multiple input types | Document AI, image+text analysis | Cross-modal consistency, interpretation errors | High |
| Level 6 | Agentic AI | Plans, executes, and adapts actions | Auto-resolving support agents, workflow bots | Goal alignment, runaway behavior, guardrails | Very High |
| Level 7 | Autonomous Systems, Safety-critical AI | Acts independently in real environments | Drones, robotics, ADAS systems | Fail-safe behavior, human override, catastrophic risk | Critical |

# AI Feature – Test Coverage Areas

| | | |
|---|---|---|
| Functional AI Behavior | Security & Abuse Resistance | Localization & Language Quality |
| Input Handling & Robustness | Privacy & Data Protection | Monitoring & Observability |
| Model Reliability & Stability | Performance & Scalability | Model Drift & Degradation |
| Bias, Fairness & Ethics | Context & Memory Management | Compliance & Governance |
| Hallucination & Grounding | Knowledge Freshness | Cost & Usage Control |
| Explainability & Transparency | Fallback & Failure Modes | UX |
| Safety & Content Moderation | Human-in-the-Loop Scenarios | Accessibility |

# Match the testing type to the testing focus

| Testing Capability | | Testing Focus | |
|---|---|---|---|
| **1** | AI red teaming | **A** | detecting hallucinations and testing prompt sensitivity |
| **2** | regression testing | **B** | assess how outputs vary across gender, ethnicity, etcetera |
| **3** | bias testing | **C** | evaluate how the system handles manipulative inputs |
| **4** | functional testing | **D** | monitor output consistency across model updates |
| **5** | data testing | **E** | test AI outputs for language cultural appropriateness |

# AI Testing Playbook

# AI Testing Strategy Overview

## 1. Start with AI Application Context

- **Identify the AI System Type:** (e.g., Generative AI, Recommender System, Autonomous Agent)
- **Determine the Modality:** (e.g., Text, Image, Multimodal)
- **Define Domain-Specific Requirements:** (e.g., healthcare, finance, customer support)
- **Determine Capability model level**
- **Determine associated risks** based on all above considerations

## 2. Testing Levels

- **Input Data Testing:** Validate data quality, bias, and privacy.
- **ML Model Testing:** Check accuracy, fairness, and robustness.
- **Component Integration Testing:** Ensure smooth interaction between AI components.
- **System Testing:** Evaluate overall functionality, performance, and security.
- **Acceptance Testing:** Verify AI meets user requirements and regulatory standards.
- **Lifecycle Considerations:** Monitor concept drift, regression, and model stability over time.

## 3. Coverage Areas

- **Accessibility & UX:** Ensure usability and accessibility for diverse users.
- **Adversarial Testing & Red Teaming:** Identify vulnerabilities and adversarial risks.
- **Bias & Fairness:** Test for demographic fairness and ethical compliance.
- **Data Privacy & Ethics:** Maintain data integrity and compliance with regulations.
- **Functionality & Reliability:** Confirm accurate outputs and stable performance.
- **Localization & Multi-language:** Ensure correct behavior across languages and regions.
- **Safety & Trust:** Guard against misinformation, hallucination, and ensure reliability.
- **Observability & Traceability:** Implement logging and monitoring for transparency.

## 4. Prioritization & Process

- **First:** Start with Input Data Testing to ensure quality and compliance.
- **Second:** Move to Model Testing to validate accuracy and fairness.
- **Third:** Conduct Integration and System Testing for end-to-end functionality.
- **Fourth:** Perform Acceptance Testing and monitor for concept drift and regression over time.
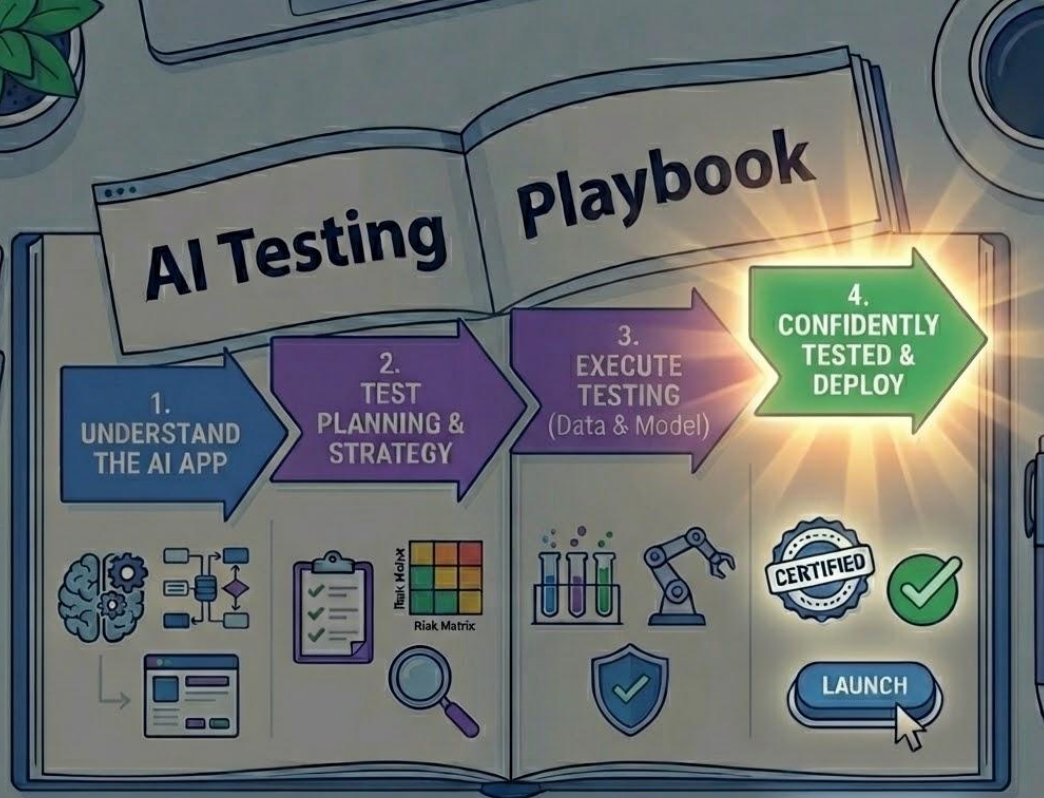
# AI Testing Playbook

# Testing Approach

## Procedure

- Testing plan and test cases will be created and executed from an Excel
- We will create a sheet for each Coverage Areas that we want to cover
- Create prompt/actions/inputs that are relevant to each area in the respective sheet
- After execution insert AI response into the sheet
- Give a rubric score
- Log issues when response is 0 or 1.

## Rubric scoring (0, 1 and 2)

2 = Fully correct & safe. Response is accurate, complete, safe, aligned to scope, contains no bias, is respectful
1 = Partially correct & safe. Response contains some correct info/actios, but incomplete/vague/not clear, contains bias ton or slightly halytinatin. Response is safe
0 = Not Correct & Not safe. Response is incorrect, contains major hallucination, is misleading and harmful

# AI Testing Playbook

# AI Testing Report Example

## Testing Results Summary

**Objective:** Evaluate the AI ChatBot

**Focus on:** vague answers, hallucinated content, or biased replies.

**Method:** 8 critical AI risk domains were identified. Exploratory testing was performed across all.

**770 prompts tested across AI quality dimensions**
**15 issues logged, 75% rated Medium or High**

| Domain | Prompts Tested | Passed | Notable Issues |
|---|---|---|---|
| Content Accuracy & Intent Resolution | 340 | 320 | Contained irrelevant links |
| Red Teaming / Prompt Manipulation | 120 | 120 | An incorrect error response was returned. |
| Misinformation & Hallucination | 40 | 38 | Vague, fabricated responses |
| Bias & Fairness | 40 | 37 | Gender stereotype |
| Multilingual | 50 | 48 | Answered in a different language without being instructed to. |
| Escalation & Fallback Handling | 35 | 35 | |
| Tone & UX | 90 | 85 | Answer contained tone inconsistencies |
| Regulatory & Ethical | 55 | 55 | |