

# AI/ML Data Readiness Checklist

# A practical guide to building reliable AI systems

Artificial intelligence is becoming part of everyday business operations, influencing analysis, reporting, and decision-making across organisations. As adoption increases, so does the need for confidence in how these systems are designed, deployed, and maintained.

This checklist has been created to help teams assess whether the foundations supporting their AI and data initiatives are ready for production use. It focuses on the practical areas that most directly influence reliability, trust, and long-term performance.

Rather than concentrating on models alone, this guide looks at the broader environment around AI, including data quality, system stability, governance, and operational readiness. These elements play a significant role in determining whether AI performs consistently once it moves beyond experimentation and into real-world use.

The checklist is designed to be:

- Clear and practical
- Relevant to both technical and business stakeholders
- Applicable across different industries and use cases
- Focused on outcomes, not theory

It can be used as a self-assessment tool, a discussion framework, or a starting point for improving existing processes. Not every organisation will be at the same stage, and not every section will apply equally. The value comes from gaining visibility into strengths, gaps, and areas for improvement.

Reliable AI is not the result of a single decision. It is built over time, through structure, clarity, and good operational practice.

This checklist is intended to support that journey.

# Data Ownership & Accountability

Understanding who is responsible for what when things break helps get services back up and running quickly. Clear ownership removes uncertainty, reduces downtime, and makes issue resolution more efficient.



**Every dataset has a named owner**

Each dataset should have a clearly identified owner who takes responsibility for its accuracy, structure, and upkeep. This creates a reliable point of contact when questions arise or changes are required.



**Data SLAs are defined (freshness, completeness, uptime)**

Setting expectations around how current, complete, and available data should be helps teams align on what “good” looks like. Clear standards reduce ambiguity and support more proactive management of data quality.



**A clear escalation process exists**

When data issues occur, it should be obvious who to contact and how to raise concerns. A well-defined escalation path improves response times and helps issues get resolved efficiently.

# Data Understanding

A strong understanding of data allows teams to work with confidence and make better decisions. When the origin, structure, and meaning of data are clear, issues are easier to identify and outcomes become more predictable.

<input type="checkbox"/>	<b>Data sources are clearly identified</b>	Each dataset should have a known origin, whether it comes from an application, external provider, or internal system. Clear source visibility makes it easier to investigate issues and improves trust in downstream analysis.
<input type="checkbox"/>	<b>Update frequency is explicitly known</b>	Teams should understand how often data is refreshed and how current it is expected to be. Clarity around freshness helps set realistic expectations and supports better decision-making.
<input type="checkbox"/>	<b>Transformation steps are understood</b>	It should be clear how data changes as it moves through systems before reaching the model. Understanding this journey makes it easier to trace issues and maintain consistency.
<input type="checkbox"/>	<b>Critical fields are clearly defined</b>	Teams should agree on which fields are essential and which are supplementary. This helps prioritise attention and ensures focus is placed on what matters most.
<input type="checkbox"/>	<b>“Bad data” has a shared definition</b>	There should be a common understanding of what poor-quality data looks like. Agreeing on this early supports faster identification and resolution of potential issues.
<input type="checkbox"/>	<b>Business meaning is documented</b>	Key fields should be explained in plain language. Clear definitions ensure that models reflect the intended business logic and reduce misunderstandings across teams.

# Data Quality Baselines

Establishing quality benchmarks makes it easier to understand when data is performing as expected and when attention is needed. These baselines provide a practical reference point for maintaining consistency over time.

<input type="checkbox"/>	<b>Expected row counts are defined</b>	Teams should have a clear view of the volume of data they normally receive. This makes unusual increases or drops easier to notice and quicker to investigate.
<input type="checkbox"/>	<b>Schema and data types are enforced</b>	Keeping data structures consistent ensures systems can process information reliably. Clear definitions reduce the risk of unexpected behaviour.
<input type="checkbox"/>	<b>Null thresholds are clear</b>	Knowing how much missing data is acceptable helps teams maintain consistent standards and focus effort where it has the greatest impact.
<input type="checkbox"/>	<b>Uniqueness rules exist</b>	Where records are expected to be unique, this should be clearly defined. Consistency here supports more accurate reporting and analysis.
<input type="checkbox"/>	<b>Valid ranges are set</b>	Defining acceptable values for numerical fields improves reliability and helps highlight unusual patterns early.

# Pipeline Resilience

Data pipelines and ML models usually run in quiet harmony. When something does go wrong, knowing where to focus attention is key to restoring services quickly and minimising disruption.

<input type="checkbox"/>	<b>Ingestion is automated</b>	Automating how data enters systems supports consistency and reduces avoidable delays or errors.
<input type="checkbox"/>	<b>Failures are immediately visible</b>	Clear visibility into issues helps teams respond quickly and reduces disruption across dependent services.
<input type="checkbox"/>	<b>Systems and dashboards provide oversight</b>	Dashboards and monitoring systems should make pipeline health easy to understand at a glance. This helps teams identify issues quickly and prioritise the right actions.
<input type="checkbox"/>	<b>Retry logic is in place</b>	Allowing systems to recover automatically from temporary issues helps maintain continuity and improves operational stability.
<input type="checkbox"/>	<b>Safe fallback behaviour exists</b>	Planning for missing or invalid data enables systems to continue operating and reduces wider impact during incidents.
<input type="checkbox"/>	<b>Latency is actively monitored</b>	Understanding where delays occur helps teams focus improvement efforts and maintain predictable performance.
<input type="checkbox"/>	<b>Upstream dependencies are known</b>	Knowing which systems influence data flow makes it easier to assess impact and coordinate resolution efficiently.

# Training vs Production Parity

Most AI failures don't come from bad models, they come from using a model on data it was never trained to understand. Consistency between training and live data is mandatory, not optional.

<input type="checkbox"/>	<b>Training and production schemas match</b>	AI and data failures are often the result of discrepancies between development, UAT and production environments. Keeping structures aligned across these systems helps reduce unexpected behaviour when models go live.
<input type="checkbox"/>	<b>Feature logic is identical</b>	Using the same transformations in training and production reduces variation and improves confidence in outputs.
<input type="checkbox"/>	<b>Data distributions are compared</b>	Monitoring how live data compares to training data helps teams recognise meaningful change and respond appropriately.
<input type="checkbox"/>	<b>Category changes are monitored</b>	Tracking new or missing values supports accurate interpretation and ongoing performance.
<input type="checkbox"/>	<b>Live data is tested before launch</b>	Testing with production-like data ahead of release helps teams validate assumptions and reduce surprises.
<input type="checkbox"/>	<b>Model assumptions are documented</b>	Recording key assumptions supports clarity, communication, and long-term maintainability.

# Dataset Versioning

Versioning helps prevent major failures when changes are introduced. When data, features, or models are updated, versioning makes it easier to understand what changed and reduces the risk of a new release breaking existing behaviour.

<input type="checkbox"/>	<b>Training data is versioned</b>	Keeping identifiable snapshots of training data allows teams to understand what a model learned from at any point in time.
<input type="checkbox"/>	<b>Feature definitions are controlled</b>	Managing changes to feature logic helps teams track how data evolves and ensures models behave as expected.
<input type="checkbox"/>	<b>Model-to-data mapping exists</b>	Linking models to the data they were trained on improves transparency and supports investigation when results change.
<input type="checkbox"/>	<b>Historical datasets are stored</b>	Retaining past versions of data supports comparison, review and analysis as systems mature.
<input type="checkbox"/>	<b>Change history is visible</b>	Clear records of what changed and when support accountability and smoother collaboration.

# Observability & Alerts

Clear visibility into systems and data helps teams act with confidence. When performance and quality are easy to see, issues can be addressed before they affect outcomes.

<input type="checkbox"/>	<b>Pipeline health is visible</b>	Teams should have a clear view of whether data is flowing as expected, delayed, or interrupted. This supports timely and informed responses.
<input type="checkbox"/>	<b>Anomalies are logged and reviewed</b>	Recording unexpected events creates a reliable history that helps teams understand patterns and improve over time.
<input type="checkbox"/>	<b>Model confidence is observed</b>	Tracking confidence alongside performance metrics provides a fuller picture of reliability
<input type="checkbox"/>	<b>Dashboards support decision-making</b>	Information should be accessible and easy to interpret. Well-designed dashboards help teams focus on what matters most.
<input type="checkbox"/>	<b>Alerts reach the right people</b>	Notifications should be directed to those best placed to respond, improving resolution times.
<input type="checkbox"/>	<b>Response steps are clear</b>	When alerts occur, teams should know what actions to take and where responsibility sits.

# Security & Compliance

Strong governance protects both your data and your reputation. Effective security and compliance processes reduce risk and support responsible use of AI across the organisation.

<input type="checkbox"/>	<b>Access is role-controlled</b>	Access to data should be based on role and responsibility. This helps ensure sensitive information is only handled by those who need it.
<input type="checkbox"/>	<b>Sensitive data is protected</b>	Personal and confidential data should be encrypted and handled appropriately. This supports trust and reduces exposure.
<input type="checkbox"/>	<b>Audit trails are in place</b>	Clear records of data access and changes support accountability and help resolve issues quickly.
<input type="checkbox"/>	<b>Data retention is defined</b>	Knowing what data to keep, and for how long, supports compliance and reduces unnecessary risk.
<input type="checkbox"/>	<b>Regulatory requirements are understood</b>	Teams should be aware of applicable obligations and reflect them in how data is handled and stored.
<input type="checkbox"/>	<b>Credentials are managed securely</b>	Passwords, keys and tokens should be stored and managed safely. This reduces the risk of accidental exposure.

# Go-Live Gate

Releasing a model into production is a significant operational step. A clear go-live process helps teams move forward with confidence and reduces the risk of disruption.

<input type="checkbox"/>	<b>Data owner approval is given</b>	Formal sign-off confirms that the data meets agreed expectations and is ready for use.
<input type="checkbox"/>	<b>Rollback is prepared</b>	Having a rollback approach in place makes it easier to restore services quickly if unexpected issues arise.
<input type="checkbox"/>	<b>Fallback behaviour exists</b>	When predictions cannot be delivered as planned, alternative processes should support continuity.
<input type="checkbox"/>	<b>Monitoring is live</b>	Performance, data quality and drift checks should be active from the outset to provide early insight.
<input type="checkbox"/>	<b>Documentation is available</b>	Clear documentation helps teams understand how systems work and how to respond when changes are needed.
<input type="checkbox"/>	<b>Ownership continues after launch</b>	Responsibility remains in place beyond release, supporting reliability as the system evolves.

Building reliable AI is not about reaching a final destination. It is about establishing practices that allow systems to evolve while remaining stable, trustworthy, and well governed.

No organisation completes this checklist perfectly on day one. Each section highlights areas where greater clarity, structure and oversight can strengthen outcomes. The aim is not perfection, but progress.

Used regularly, this checklist can support:

- Better decision-making
- Faster issue resolution
- Stronger collaboration between business and technical teams
- Greater confidence in AI-driven outcomes

As AI systems become more embedded in everyday operations, the foundations that support them make an increasing difference to long-term success.

## **About Shipshape Data**

Shipshape Data supports organisations in building dependable data foundations for AI and analytics. We work with teams to improve data quality, strengthen pipelines, and increase operational visibility, helping AI systems perform more consistently in real-world environments.

If you would like to discuss any areas highlighted in this checklist, you can find more information about our services at:

[www.shipshapedata.com](http://www.shipshapedata.com)