

✦ WHITE PAPER

MEASURING THE TRUE PERFORMANCE OF AI AGENTS

A data-driven framework for evaluating AI agent reliability, user experience, and real-time CX impact.

Ralf Ellspermann, CSO
2025



Executive Summary

The rapid adoption of AI agents in customer experience (CX) has created a significant measurement gap. Traditional contact center metrics, while still relevant, are insufficient for evaluating the nuanced performance of AI-driven systems. As a result, fewer than 20% of organizations currently track well-defined KPIs for their generative AI solutions [1]. This lack of visibility creates a significant risk, as organizations are flying blind, unable to distinguish between true value creation and the illusion of progress. The service organization of the future is a partnership between humans and machines, and this partnership needs to be evaluated on how well it works together, not just how fast it moves.

This document presents a new framework for measuring the true performance of AI agents, moving beyond simplistic accuracy scores to a more holistic and actionable understanding of their impact. We introduce the two pillars of AI agent measurement: empirical metrics (the “what”) and experiential metrics (the “how”). We then provide a deep dive into a comprehensive set of AI-native and reliability metrics, including containment, escalation quality, and consistency. Finally, we outline a practical approach for building real-time CX dashboards based on a three-tier governance framework, enabling organizations to protect their brand, drive operational value, and scale with confidence.

By adopting this new measurement framework, organizations can move from managing activities to orchestrating outcomes, fostering a culture of continuous improvement, and unlocking the full potential of their AI investments. This is not just about better dashboards; it’s about building a more resilient, efficient, and customer-centric service organization for the future.

The Two Pillars of AI Agent Measurement

To effectively measure the performance of AI agents, a balanced approach is required, one that considers both the objective, system-level data and the subjective, user-level feedback. These two pillars, empirical and experiential metrics, provide a comprehensive view of AI agent

performance, ensuring that organizations are not just measuring what is easy, but what is important.

Empirical Metrics: The “What”

Empirical metrics are the objective, quantifiable data points that describe what the AI system is doing. They are the foundation of any robust measurement framework, providing a clear and unbiased view of the agent’s performance. These metrics include:

- **Accuracy:** The correctness of the agent’s responses and actions.
- **Resolution Rate:** The percentage of interactions successfully resolved by the agent without human intervention.
- **Escalation Rate:** The frequency with which the agent escalates interactions to a human agent.
- **Latency:** The time it takes for the agent to respond to user inputs.
- **Error Rate:** The frequency of incorrect intents, misrouted flows, or wrong data returned.

These metrics are essential for understanding the core functionality of the AI agent and identifying areas for improvement. They provide the hard data needed to make informed decisions about system configuration, training data, and prompt engineering.

Experiential Metrics: The “How”

Experiential metrics, on the other hand, capture the subjective experience of the user. They provide insight into how it feels to interact with the AI agent, which is often a leading indicator of customer satisfaction and loyalty. These metrics include:

- **Clarity:** The ease with which users can understand the agent’s responses.
- **Effort:** The amount of work required from the user to resolve their issue.
- **Trust:** The user’s confidence in the agent’s ability to provide accurate and reliable information.
- **Satisfaction:** The user’s overall satisfaction with the interaction.
- **Return Usage:** The likelihood that the user will interact with the agent again in the future.

Experiential metrics are typically gathered through post-interaction surveys, sentiment analysis, and behavioral proxies (e.g., rephrasing, repeated queries). They provide the qualitative context needed to understand the “why” behind the empirical data.

The Balance: A Holistic View of Performance

As ASAPP aptly puts it, “To effectively measure your generative AI agent, you need both types of metrics. Measuring only one is how failed pilots go undiagnosed until you hit scale and start losing customers” [2]. An AI agent can have a high resolution rate but a low satisfaction score, indicating that while it is technically resolving issues, it is doing so in a way that is frustrating or confusing for users. Conversely, an agent can have a high satisfaction score but a low resolution rate, suggesting that while users enjoy interacting with it, it is not actually solving their problems.

By combining empirical and experiential metrics, organizations can gain a holistic view of AI agent performance, enabling them to identify and address issues before they impact the customer experience. This balanced approach is the cornerstone of a successful AI measurement strategy, providing the insights needed to build a truly effective and customer-centric AI-powered service organization.

A Deep Dive into AI-Native Metrics

While traditional contact center metrics provide a useful starting point, the unique capabilities and challenges of AI agents require a new set of AI-native metrics. These metrics are designed to measure the specific aspects of AI performance that have the greatest impact on the customer experience and operational efficiency.

Core Empirical Metrics

First Contact Resolution (FCR)

FCR remains the gold standard for measuring the effectiveness of any customer service interaction, and it is particularly important for AI agents. A high FCR rate indicates that the AI agent is able to successfully resolve customer issues without the need for human intervention, which is the ultimate goal of any automation initiative. As ASAPP notes, “It’s the most honest proxy for whether your generative AI agent works” [2].

Error Rate

Uncaught errors can have a significant impact on the customer experience and can quickly erode trust in the AI agent. It is essential to track a variety of error types, including incorrect intents, misrouted flows, and wrong data returned. By identifying and addressing these errors in a timely manner, organizations can prevent them from compounding and causing more significant issues down the line.

Containment Rate

Containment rate measures the percentage of interactions that are fully handled by the AI agent without being escalated to a human agent. While a high containment rate is generally desirable, it is important to guard against over-containment at the cost of the customer experience. As ASAPP warns, “Abandonment and escalation are signals of failure” [2].

Escalation Frequency & Quality

When an AI agent does need to escalate an interaction to a human agent, it is important that the handoff is as seamless as possible. This requires not only tracking the frequency of escalations, but also the quality of the escalation. A high-quality escalation includes all of the relevant context from the AI interaction, enabling the human agent to quickly understand the issue and provide a resolution.

Latency & Response Time

In the age of instant gratification, latency can have a significant impact on the customer experience. As ASAPP points out, “Delay is the enemy of confidence, especially in voice interactions” [2]. It is important to track both the time to first response and the time to resolution, with the goal of providing a consistently responsive experience.

Key Experiential Metrics

CSAT/NPS

Customer satisfaction (CSAT) and Net Promoter Score (NPS) are two of the most widely used experiential metrics, and they are just as relevant for AI agents as they are for human agents. By tracking CSAT and NPS for AI-led interactions, organizations can gain valuable insights into how customers perceive the AI agent and identify areas for improvement.

Customer Effort Score (CES)

CES measures the amount of effort required from the customer to resolve their issue. A low CES score indicates that the AI agent is easy to use and provides a frictionless experience. As with CSAT and NPS, it is important to track CES for AI-led interactions to ensure that the agent is not creating unnecessary friction for customers.

Trust & Confidence Signals

Trust is a critical component of any successful AI implementation. If customers do not trust the AI agent, they are unlikely to use it, regardless of its technical capabilities. Trust can be measured through a variety of signals, including drop-offs after vague messages, rephrasing of questions, and repeated queries. By tracking these signals, organizations can identify and address issues that may be eroding customer trust.

Reliability Metrics: Moving Beyond Accuracy

While accuracy is a fundamental metric for evaluating AI agent performance, it is not sufficient on its own. A high-accuracy agent can still fail in production if it is not reliable. As Galileo notes, “An agent might achieve 95% accuracy on your test set but still fail catastrophically in production due to poor consistency, inability to handle edge cases, or performance degradation over time” [3]. To ensure that AI agents are truly production-ready, organizations must also measure a variety of reliability metrics.

Consistency and Determinism

Consistency is the ability of an AI agent to provide similar responses to similar questions, even when the wording is different. This is particularly challenging for LLMs, which are inherently non-deterministic. As Charity Majors, CTO of Honeycomb, has noted, “traditional computers are really good at very precise things and very bad at fuzzy things, our LLMs are, like, really bad at pretty very precise things and really good at fuzzy things” [3]. To address this challenge, organizations must create systematic testing approaches that simulate how real users communicate, rather than how developers think they should.

Robustness Under Adversarial Conditions

Robustness is the ability of an AI agent to handle unexpected or malformed inputs gracefully. This includes everything from typos and grammatical errors to more malicious attempts to break the system. As Galileo points out, “These aren’t edge cases; they happen from time to time in

any production environment” [3]. To build robust AI agents, organizations must deliberately break things in controlled ways, creating scenarios that mirror how systems fail in production.

Uncertainty Quantification & Confidence Calibration

A well-calibrated AI agent knows when it doesn’t know the answer and is able to communicate its uncertainty to the user. This is critical for building trust and preventing the agent from providing confident-sounding but incorrect information. As Galileo explains, “When your agent says it’s 90% confident, it should be correct about 90% of the time. When it expresses uncertainty, that uncertainty should correlate with situations where human oversight helps” [3].

Temporal Stability & Performance Drift

Performance drift is the gradual degradation of an AI agent’s performance over time. This can be caused by a variety of factors, including changes in user behavior, data distributions, and infrastructure. As Galileo warns, “This is the slow degradation that represents one of the most insidious threats to production AI systems” [3]. To combat performance drift, organizations must establish baseline measurements across key reliability indicators and track their evolution over time.

Context Retention & Coherence

Context retention is the ability of an AI agent to maintain relevant information across extended conversations. This is essential for creating a natural and purposeful user experience. As Galileo notes, “Nothing frustrates users more than agents that seem to forget important conversation details or contradict themselves within the same interaction” [3].

Response Latency Consistency

While speed is important, predictability is often more so. As Galileo explains, “Users can adapt to consistently slower responses, but they struggle with systems that respond instantly sometimes and take forever other times without any apparent reason” [3]. By focusing on reducing variability within acceptable ranges, organizations can create a more predictable and trustworthy user experience.

Graceful Degradation Under Load

Graceful degradation is the ability of an AI agent to maintain core functionality during periods of high demand or infrastructure issues. As Galileo notes, “The difference between systems that fail gracefully and those that fail catastrophically becomes critical during periods of peak demand” [3].

Behavioral Consistency Across Demographics

Finally, it is essential to ensure that AI agents provide a consistent experience for all users, regardless of their communication style, cultural background, or accessibility needs. This requires a commitment to fairness and equity in all aspects of AI development and deployment.

Building Real-Time CX Dashboards

To effectively manage the performance of AI agents, organizations need more than just a list of metrics; they need a comprehensive and real-time view of their entire CX ecosystem. This requires a new generation of CX dashboards that are designed to provide a holistic view of the human-machine partnership.

A Three-Tier Governance Framework

A successful CX dashboard should be built around a three-tier governance framework that aligns with the organization’s strategic priorities:

- **Tier 1: Protect the Brand:** This foundational tier focuses on the metrics that are most critical for protecting the brand and ensuring a positive customer experience. These include containment, error rate, and latency.
- **Tier 2: Drive Operational Value:** This tier focuses on the metrics that are most important for driving operational efficiency and ROI. These include CSAT, effort, and escalation frequency.
- **Tier 3: Scale with Confidence:** This tier focuses on the metrics that are most critical for scaling the AI solution and driving continuous improvement. These include learning velocity, retention, and trust signals.

The Technology Stack for Observability

Building a real-time CX dashboard requires a modern technology stack that is designed for observability. This includes:

- **Data Collection & Instrumentation:** The foundation of any dashboard is a robust data collection and instrumentation strategy. This requires logging every interaction, decision point, and user action to create a rich and detailed dataset.
- **Analytics & Visualization:** The next layer is a powerful analytics and visualization engine that can transform raw data into actionable insights. This includes real-time dashboards, alerting, and trend analysis.
- **Feedback Loops:** Finally, it is essential to create a tight feedback loop between the AI agent and the human agents who are using it. This enables human agents to provide real-time feedback on the AI's performance, which can be used to drive continuous improvement.

The Future of Service Measurement

The age of AI requires a new approach to service measurement. Traditional metrics are no longer sufficient to capture the complexity and nuance of the modern, AI-powered contact center. By embracing a new framework based on AI-native and reliability metrics, organizations can gain a deeper and more holistic understanding of their AI agent's performance, enabling them to move beyond managing activities to orchestrating outcomes.

This is not just about better dashboards; it is about building a more resilient, efficient, and customer-centric service organization. It is about fostering a culture of continuous improvement, where data is used to drive every decision. And it is about recognizing that the future of customer service is not about replacing humans with machines, but about creating a powerful partnership between the two.

Contact Ralf Ellspermann, CSO, to discuss how your organization can build a unified performance governance model for AI-powered service operations and ensure measurable ROI across human–AI collaboration.

References

[1] Oracle. (2023). Measuring What Matters in the Age of AI Agents.

[2] ASAPP. (2023). How to measure your generative AI agent performance (and why you can't afford to get this wrong).

[3] Galileo. (2023). 8 AI Agent Metrics That Go Beyond Accuracy.

© 2025 Ralf Ellspermann, CSO, CynergyCx.ai.

This white paper contains proprietary research, performance frameworks, and analytical insights developed through extensive industry expertise.

Reproduction or distribution without permission is strictly prohibited.