

Generative AI ROI and Productivity Assessment

Presentation for Mastercard AI Governance Unit

04/28/2025

Ethan Barrilleaux, Jada Battle, Lambert Li, Makeenie Robinson, Aidan Whitlock

Wake Forest MSBA Practicum Team

MEET OUR TEAM



Makeenie Robinson

Team Lead



Jada Battle

Client Liaison



Aidan Whitlock

Advisor Liaison



Ethan Barrilleaux

Design Lead



Lambert Li

Research and Knowledge Officer

AGENDA

01

Context and Insights

- Gaps in current scorecard for GenAI ROI
- Industry flags: ROI missteps & risk blind spots
- Need for business-aligned metrics

02

Our Solution

- User trust + model performance metrics
- Logic-embedded questions, audience clarity
- Anchoring fairness, efficacy, transparency

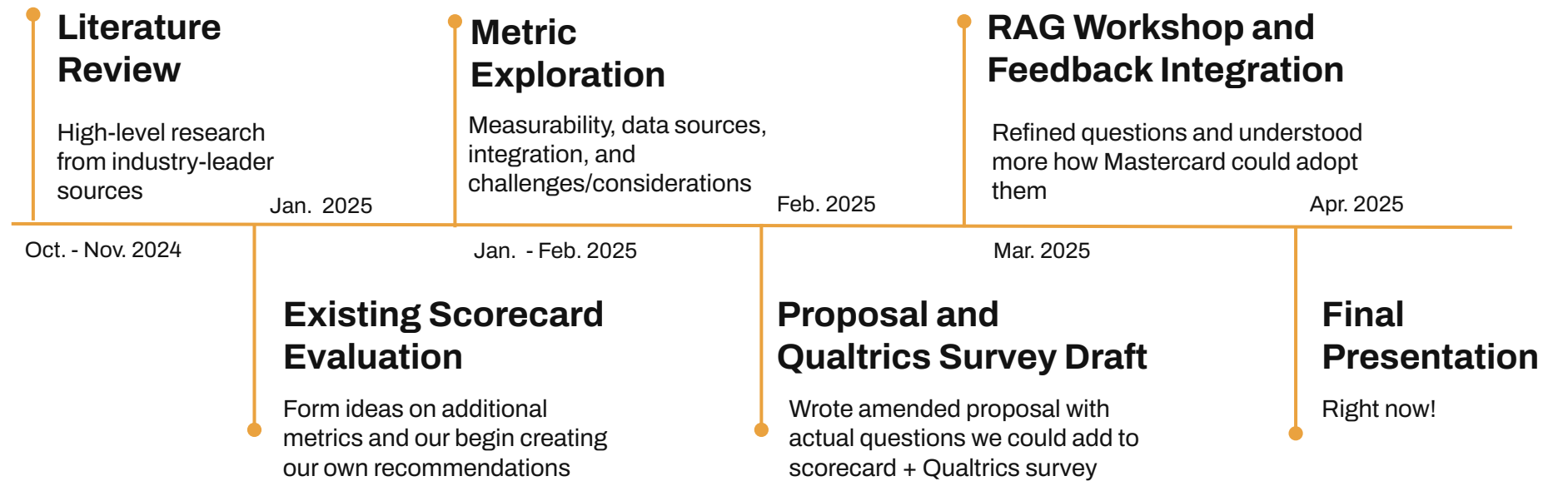
03

Final Deliverables

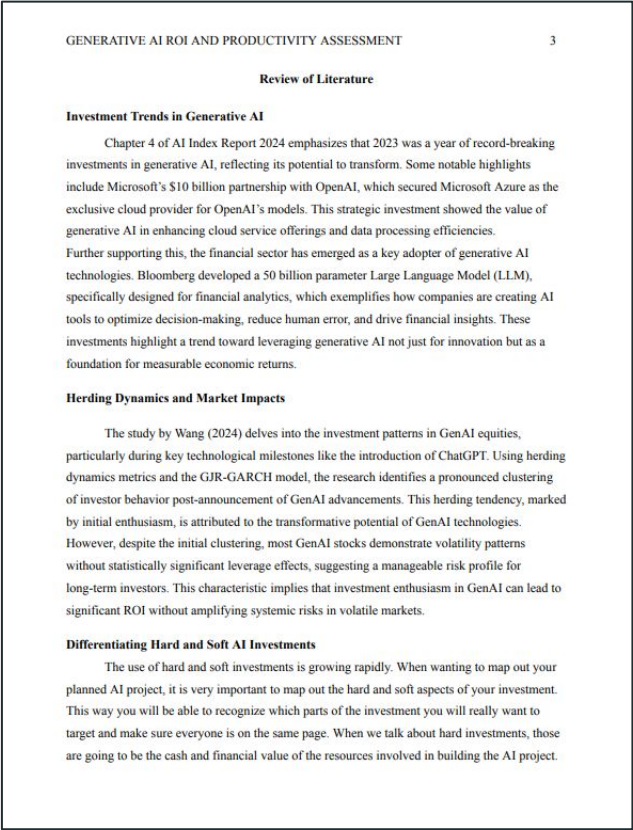
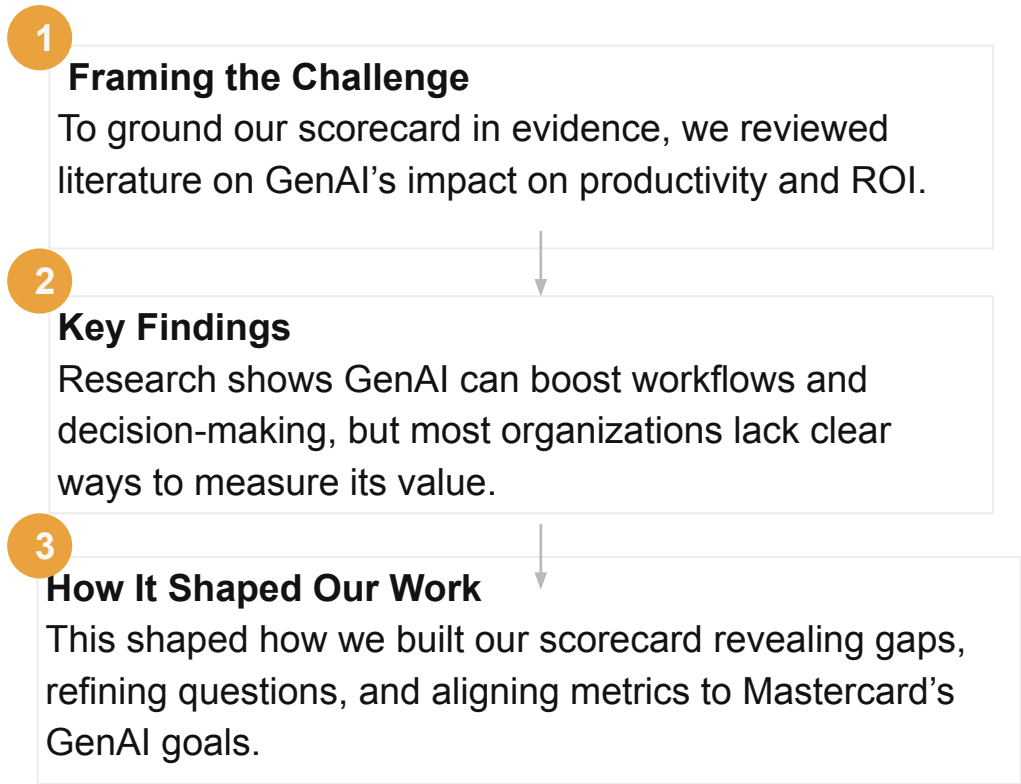
- Actionable, logic-driven question set
- Excel tracker + Qualtrics survey demo
- Clear guidance for internal + vendor rollout



From Kickoff to Delivery: Our End-to-End Scorecard Development Timeline (Oct. 2024–Apr. 2025)



An overview from the literature on Generative AI’s business value and productivity impact.



Analyzing the existing scorecard to identify gaps and inform targeted improvements.

- ✓ We identified gaps in Mastercard’s risk scorecard around GenAI. Our solution: tailored, clear, and targeted questions that better align with GenAI risks and business needs.

Existing Approach	Identified Gaps	Our Approach
Post-build evaluation focus	Misses early stage risk planning and alignment	Integrated Pre/Post build tags to clarify question timing and relevance
Broad, AI generic questions	Fails to capture the unique challenges of GenAI	Developed GenAI specific metrics
Technical and complex phrasing	Difficult for business owners to interpret and act upon	Applied audience mapping and simplified language for clarity



Understand how metric categories drive design philosophy.

Why Separate Them?

✓ A GenAI system may perform very well technically, but have poor ROI if employees do not use, understand, or trust the model.

Recommendation: Implement a Decision Confidence Score (DCS) to measure confidence in AI-generated recommendations. Integrate DCS tracking into governance pipeline through post-interaction surveys and override monitoring.

User Trust & Risk Assessment:

- How confident are you in the AI-generated recommendation? (Scale: Not at all – Completely confident)
- Have you overridden an AI recommendation due to a lack of trust? (Yes/No)

1.5 Friction Index (FI) Integration

Recommendation: Implement the Friction Index (FI) to quantify inefficiencies in AI interactions, ensuring smoother and more efficient user experiences.

Proposal: Track user journey data, including task completion times, abandonment rates, and feedback loops to identify areas of friction in AI-driven systems.

User Interaction & Friction Assessment:

- How frequently do you encounter difficulties or delays when interacting with AI-powered systems? (Scale: Never – Very Often)

2. Model Metrics

2.1 Perplexity & Long Perplexity (PPL) Management

Recommendation: Monitor perplexity scores to detect model inconsistencies and generate coherent AI-generated text.

Model Perplexity Evaluation: How well does your AI model maintain coherence in long-form text generation? (Scale: Poor - Excellent)

- Please provide the most recent Perplexity Score and Long Perplexity Score.

2.2 Prediction Entropy Reduction

User Interaction / Satisfaction Metrics

1.1 Hallucination Rate Reduction

Recommendation: Implement Retrieval-Augmented Generation (RAG) and confidence scoring to minimize hallucinated AI responses.

(Prototype question for Qualtrics)

Hallucination Rate Control: How frequently does your AI model generate incorrect or misleading outputs? (Scale: Never - Very Often)

- Please provide the latest measured Hallucination Rate. See a sample from the appendix.

1.2 User Override Rate Optimization

Recommendation: Use SHAP-based explainability tools to improve AI recommendation acceptance. Assess AI-response quality from User Override Rate (if any).

How do you currently assess user confidence in AI-driven decisions? (Multiple Choice: Post-interaction surveys, Override monitoring, Qualitative feedback, No structured assessment, Other)

User Override Rate Assessment: How often do human reviewers override AI-generated recommendations? (Scale: Never - Very Often)

- Please provide the User Override Rate percentage. See a sample from appendix.

1.3 AI-Driven Value Attribution Score Tracking

Recommendation: Measure the impact of AI-driven decisions on business outcomes such as cost savings and revenue growth.

Proposal: Deploy A/B testing frameworks to track and compare AI-driven vs. human-driven results.

Business Impact Validation: Can you provide quantifiable evidence of AI-driven improvements in business performance? (Yes/No)

- If Yes, please provide supporting metrics (e.g., revenue impact, cost reduction, user engagement improvement).

1.4 Decision Confidence Score (DCS) Implementation

Incorporate metrics to better quantitative measurements.

1 What Did We Do?

We designed a scorecard and survey by identifying key AI metrics, mapping them to specific questions, and organizing everything into a master Excel document for streamlined analysis.


Ques- tion ID	Subject	Category	Answer Format	Aud- ience	Intent	Data Type Needed	So What?	Implica- tion
1	Hallucination Rate	User Interaction/ Satisfaction Metric	Scale	Business Owner/ Developer	To assess the factual reliability of GenAI outputs	Model outputs, benchmark facts, user feedback	High hallucination = trust/risk concern	High hallucination rate implies misinformation, reputational risk, and low adoption
2	User Override Rate	User Interaction/ Satisfaction Metric	Scale	Product Manager/ Analyst	Track user trust & alignment between AI suggestions and human judgment	System logs, user flags, override actions	High override = low trust	Frequent overrides may signal explorability or accuracy issues
3	AI Driven Value Attribution	User Interaction/ Satisfaction Metric	Binary	Business Owner/ Developer	Link AI use to measurable ROI	Business KPIs, sales data, operational cost data	Supports investment decisions	ROI linkage enables prioritization of AI projects
4	Friction Index	User Interaction/ Satisfaction Metric	Scale	Product Manager/ Developer	Identify inefficiencies in AI interactions	User journey data, task completion times, error logs	High friction = user frustration and inefficiency	High friction access indicates usability barriers that need optimization
5	Decision Confidence Score	User Interaction/ Satisfaction Metric	Scale	Business Owner/ Analyst	Assess user trust in AI decisions	Post-interaction surveys, overall monitoring	Low confidence = trust/risk concern	Lack of confidence may indicate issues with AI explainability or accuracy
6	Propriety Score	Model Performance Metric	Scale	Developer or AI Auditor	Measure text coherence among model iterations	Text-based model output, benchmark (quality standard)	Indicates when LLM may start to deviate from coherent text generation	Lower propriety score = more coherent and stable AI generated text
7	Response Consistency Score	Model Performance Metric	Numeric	Developer or AI Auditor	Quantify consistency performance	Response tracking logs, audit reports	Helps validate response reliability	Provides benchmark for consistency tuning and monitoring
8	Divergence Score	Model Performance Metric	Scale/Text	Developer or AI Auditor	Measure content divergence across model iterations	Generated text from LLM and industry standard reference text	Implement automated drift detection for model output	Lower KI and JS Divergence score = less drift

2 From Metrics to Questions

Question Identifier	Subject	Category	Question Body	Potential Answers	Answer Format	Audience	Pre/ Post?	Required to answer?	Condition	Preceding question #	Next question #	Intent	Data Type Needed	So What?	Implication
M1.a	Hallucination Rate	User Interaction/ Satisfaction Metric	How well does your AI model at avoiding incorrect or misleading outputs?	Extremely incompetent to extremely competent (7 point scale)	Multiple Choice	Business Owner/ Developer	Post	Yes	If yes in Q33	Q6b.What potential metrics do you intend to use to track success?		To assess the factual reliability of GenAI outputs	Model outputs, benchmark facts, user feedback	High hallucination = trust/risk concern	High hallucination rate implies misinformation, reputational risk, and low adoption

Sample Implementation: Bring The Scorecard To Life In Qualtrics

Restart Survey Place Bookmark Tools Share Preview

 WAKE FOREST UNIVERSITY

MI.a How well or not does your AI model at avoiding incorrect or misleading outputs?

☐ Extremely incompetent

☐ Moderately incompetent

☐ Slightly incompetent

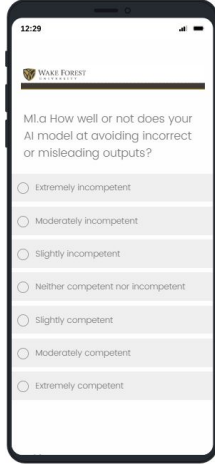
☐ Neither competent nor incompetent

☐ Slightly competent

☐ Moderately competent

☐ Extremely competent

MI.b
Please provide the latest measured Hallucination Rate.
See a sample code to measure
here: <https://drive.google.com/file/d/1ffq3MpAgIFpJLeruzY9x5Nhkw4fSIhs7/view?usp=sharing>



https://wakeforest.qualtrics.com/jfe/form/SV_cvI2NGCTVMAXvwi

Visualize how questions will be added to Mastercard framework.

1 Clear Placement w.r.t Scorecard

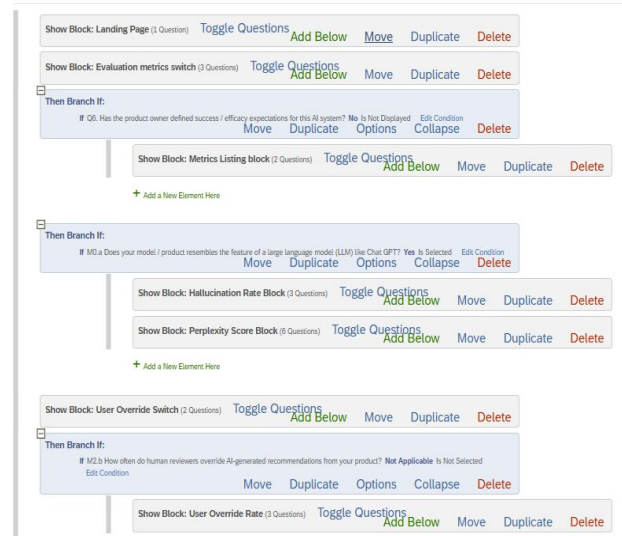
Identify existing scorecard area that have attempted to capture quantitative metrics and supplement with additional measurements in:

- User-centric / end product metrics
- Model performance metrics

Provide added clarity on how proposed metrics and merge seamlessly with existing system.

2 Conditional Logic

Survey flow Published



Incorporate conditional logic to better target applicable survey users to dedicated questions:

- Ensure accurate metrics being highlighted and areas of interested focused.
- Raise awareness on business metrics in early product dev.
- Prevent irrelevant and inaccurate information gathering.

Visualize how questions will be added to Mastercard framework.

3 Implementation Sample

M3.d If you have supporting data, please provide the friction index by referencing the following
sample: https://drive.google.com/file/d/10QZGk9kHshP_7W00FyB4X4vdRILTAH93/view?usp=sharing

```
class FrictionIndexMonitor:
    def __init__(self):
        self.error_counts = [] # Tracks AI recommendation errors
        self.resolution_times = [] # Time taken to resolve AI-related issues
        self.drop_off_rates = [] # Tracks user drop-off rates due to AI inefficiencies

    def record_interaction(self, errors, resolution_time, drop_off_rate):
        """Records an AI interaction event for friction index calculation."""
        self.error_counts.append(errors)
        self.resolution_times.append(resolution_time)
        self.drop_off_rates.append(drop_off_rate)

    def calculate_friction_index(self):
        """Computes the Friction Index (FI) as a weighted score."""
        if not self.error_counts:
            return None

        avg_errors = sum(self.error_counts) / len(self.error_counts)
        avg_resolution_time = sum(self.resolution_times) / len(self.resolution_times)
        avg_drop_off = sum(self.drop_off_rates) / len(self.drop_off_rates)

        # Assigning weights to factors (customizable based on impact analysis)
        fi_score = (0.4 * avg_errors) + (0.3 * avg_resolution_time) + (0.3 * avg_drop_off)
        return fi_score
```

Create sample implementation code to run business-centered test to help complete the scorecard.

4 Interpretable Metrics

So What?	Implication
High hallucination = trust/risk concern	High hallucination rate implies misinformation, reputational risk, and low adoption
High override = low trust	Frequent overrides may signal explainability or accuracy issues
Supports investment decisions	ROI linkage enables prioritization of AI projects
High friction = user frustration and inefficiency	High friction scores indicate usability barriers that need optimization

Each question is followed by a note of metrics result. More explorations can be done for setting up successful benchmarks from past projects.

Currently, no industry standard on particular scores.

Final Outputs Designed to Support Deployment and Decision-Making

- ✓ Master Excel Sheet of 20+ Tagged Questions
- ✓ Built-in Conditional Logic
- ✓ Demo Survey in Qualtrics
- ✓ Documents With Metrics, Use Cases, and Definitions

User Interaction / Satisfaction Metrics

1.1 Hallucination Rate Reduction

Recommendation: Implement Retrieval-Augmented Generation (RAG) and confidence scoring to minimize hallucinated AI responses.

(Prototype question for Qualtrics)

Hallucination Rate Control: How frequently does your AI model generate incorrect or misleading outputs? (Scale: Never - Very Often)

○ Please provide the latest measured Hallucination Rate. See a sample from the appendix.

1.2 User Override Rate Optimization

Recommendation: Use SHAP-based explainability tools to improve AI recommendation acceptance. Assess AI response quality from User Override Rate (if any).

By assess user confidence in AI-driven decisions? (Multiple Choice: Yes, Override monitoring, Qualitative feedback, No structured

Assessment: How often do human reviewers override AI-generated (Scale: Never - Very Often)

is the User Override Rate percentage. See a sample from appendix.

Friction Score Tracking

(the impact of AI-driven decisions on business outcomes such as cost

g. Turnovers to track and compare AI-driven vs. human-driven

Idation: Can you provide quantifiable evidence of AI-driven

since performance? (Yes/No)

provide supporting metrics (e.g., revenue impact, cost reduction, most improvement).

Score (DCS) Implementation

WAKE FOREST UNIVERSITY

The following question will be followed by Q6b or ideally replace Q6b

Q6. Has the product owner defined success / efficacy expectations for this AI system?

☐ Yes

☐ In draft

☐ No

M0.a Does your model / product resembles the feature of a large language model (LLM) like Chat GPT?

☐ Yes

☐ No

Question	Subject	Category	Answer	Audience	Intent	Data Type	So What?	Implication
2	User Override Rate	User Interaction/	Scale	Product	Track user	System logs, user	High override =	Frequent
3	AI Driven Value Attribution	User Interaction/	Binary	Business	Link AI use to	Business KPIs,	Supports	ROI linkage
4	Friction Index	User Interaction/	Scale	Product Manager/	Identify inefficiencies in completion times, error logs	User journey data, task completion times, error logs	High friction = user frustration and inefficiency	High friction barriers that need optimization. Lack of confidence may indicate issues with AI explainability or accuracy. Lower. Provides benchmark for consistency training and monitoring. Lower KL and
5	Decision Confidence Score	User Interaction/	Scale	Business Owner/	Assess user trust in AI decisions	Post-interaction surveys, override monitoring	Low confidence = trust/risk concern	Indicates when
6	Perplexity Score	Model Performance	Scale	Analyst/	Develop	Measure text	Test-based	Indicates when
7	Response Consistency Score	Model Performance	Numeric	AI Auditor	Develop	Quantify consistency	Response tracking logs, audit reports	Helps validate response reliability
8	Divergence Score	Model Performance	Scale/Text	AI Auditor	Develop	Measure	Generated text	Implement
9								
10								
11								
12								
13								

Ready Accessibility: Good to go



Appendix

1. Literature Review: [Literature Review](#)
2. Metrics Written Report: [Metrics Written Report](#)
3. Amended Proposal: [Amended Proposal](#)
4. Question Design (Qualtrics Survey): [Qualtrics Survey](#)
5. Final Excel Spreadsheet: [Final Deliverable](#)

