

77% Expects Generative AI to have the largest impact on their business out of all emerging technologies

71% Will implement their first Generative AI solution within next two years

73% Believe Generative AI will increase workforce productivity

64% Believe Generative AI will help their business gain a competitive advantage over competitors

Exhibit 6: Top priority functions for adoption are IT/Tech and operations



Notes: Sum will not add up to 100% as it is a multi-select question.

Sources: Generative AI survey, March 2023.

SRE 2.0 : Amplifying Reliability with GenAI

Indika Wimalasuriya

SRE - 2024



Agenda

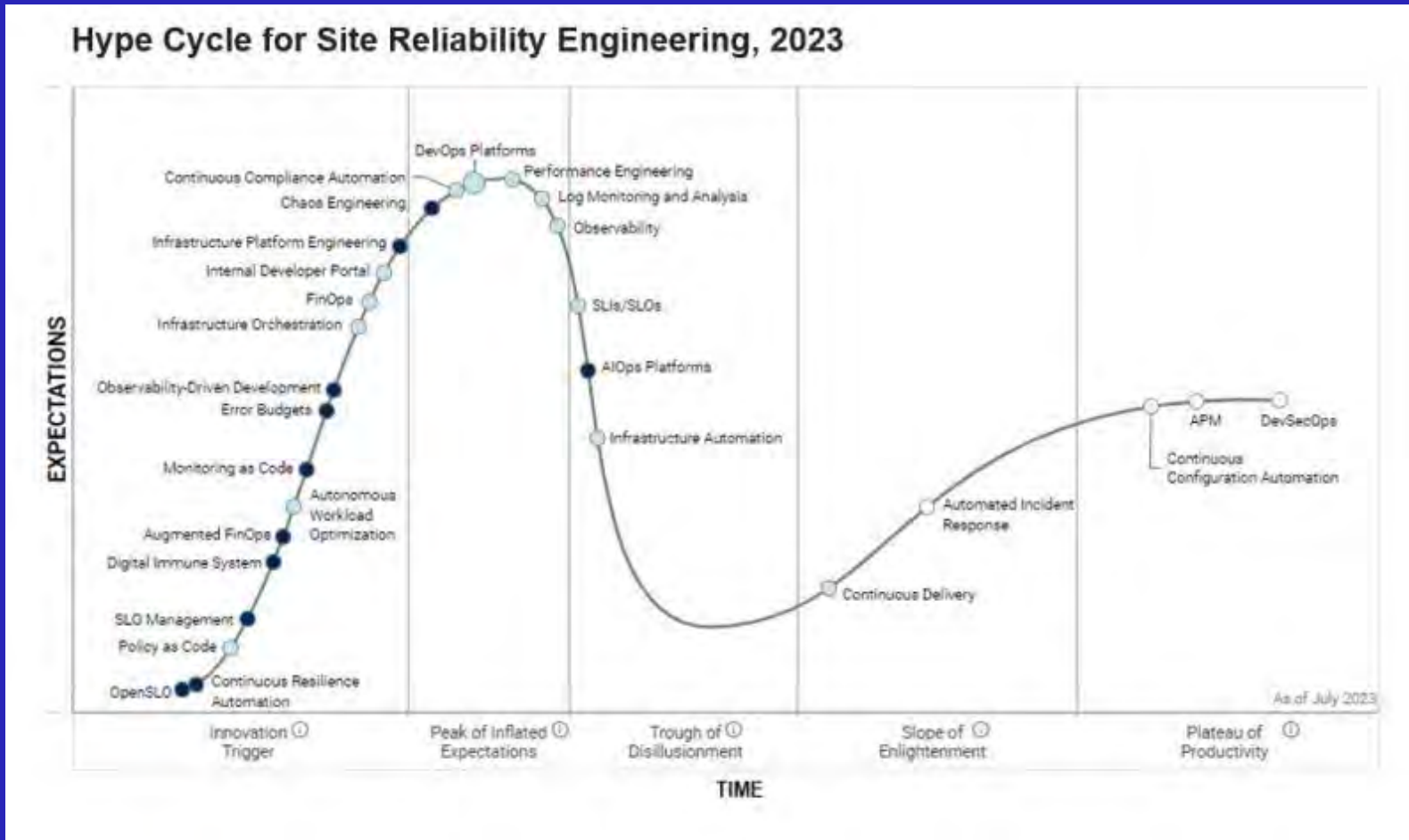
- Challenges in Traditional SRE.
- The Role of GenAI in SRE
- GenAI Impacts in Key Pillars of SRE
- GenAI Use Cases and Potential Benefits
- Implementation Strategies
- Best Practices and Pitfalls to Avoid
- Future Outlook

Quick Intro about myself



- **Resides in Colombo, Sri Lanka, with my beautiful daughter and wife.**
- **Reliability Engineering Advocate, Solution Architect (specializing in SRE, Observability, AIOps, & GenAI).**
- **Employed at Virtusa, overseeing technical delivery and capability development.**
- **Passionate Technical Trainer.**
- **Energetic Technical Blogger.**
- **AWS Community Builder - Cloud Operations.**
- **Ambassador at DevOps Institute (PeopleCert).**

Gartner SRE Hype Cycle



- Focus on trajectory of emerging technologies and frameworks for SRE practices
- Aims to ensure reliability, maximize performance, and minimize resources

Site Reliability Engineering

Key pillars of Success

- Reduce organizational Silos
- Accept failures as normal
- Implement gradual changes
- Leverage tooling and automating
- Measure everything

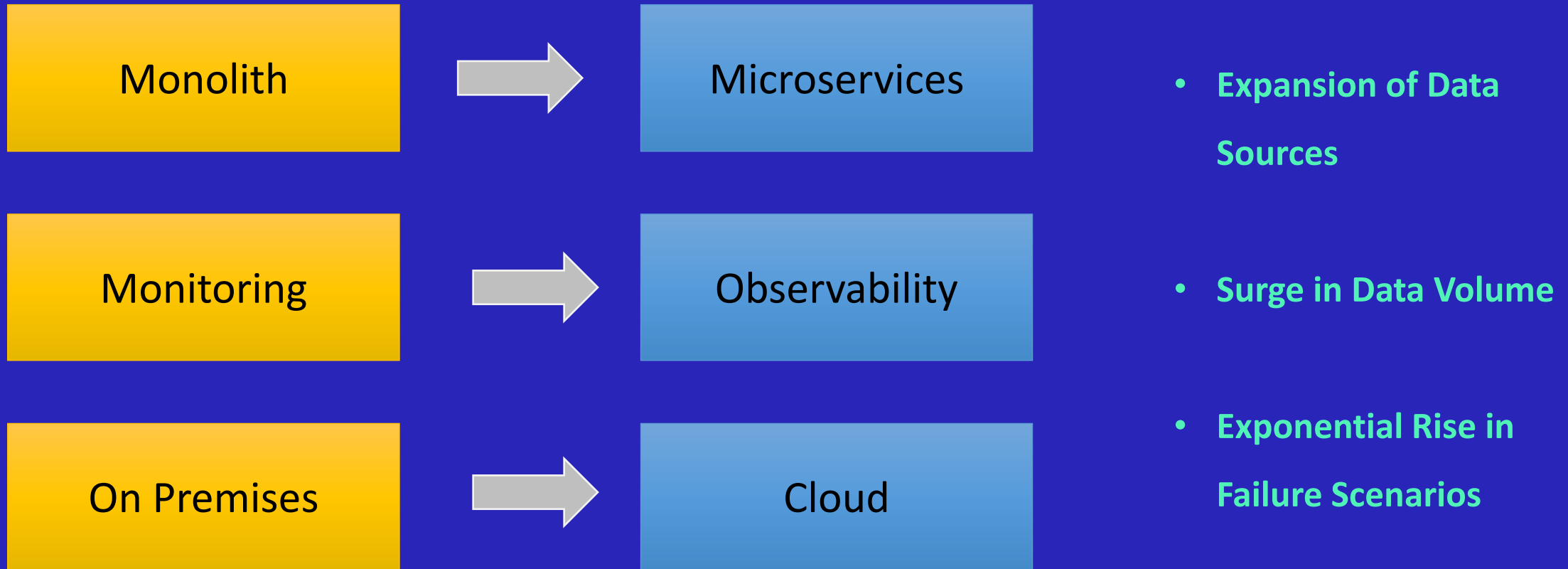
Principles

- Observability
- SLI, SLO and Error Budgets
- System architecture and Recovery

Objectives

- Release & Incident Engineering
- Automation
- Resilience Engineering
- Blameless Postmortems

Navigating Digital Transformation: Managing Ever-Growing Complexity



Operations is a Software Problem

By 2026,

Code written by developers/humans will be reduced by 50% due to generative AI code generation models

GenAI Emerges: Unveiling the Power of Next-Gen Artificial Intelligence

Generative AI, refers to machine learning models that can create new, original content like text, images, audio and more



Text Generation:

- Articles & stories - GPT-3 can generate news articles and fiction stories
- Code - GitHub Copilot suggests context-relevant code for developers



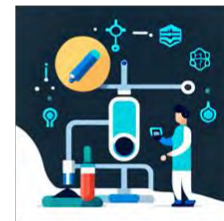
Image Generation:

- Art - DALL-E 2 creates original digital art from text prompts
- Photos - Generative models can edit or enhance photos



Video Generation:

- AI models can generate or edit video content



Drug Discovery:

- Models can analyze chemical compounds and suggest new drug candidates



Other Creative Applications:

- Generating logos, recipes, fashion designs, architectural drawings
- Personalizing content like customized ads or product recommendations

Unveiling the Potential: The Capabilities of LLM (Large Language Models)

What can LLMs do?

Inputs

Natural Language
Structured Data
Multi-Lingual text
Transcription
Computer Code

Operations

Text or Code Generation
Text Completion
Text Classification
Text Summarization
Text Translation
Sentiment Analysis
Text Correction
Text Manipulation
Named Entity Recognition
Question Answering
Style Translation
Format Translation
Simple Analytics

Outputs

Natural Language Text
Structured Data
Multi-Lingual Text
Computer Code

Navigating Challenges: Risks Associated with Large Language Models

Model Risks

- Bias
- Misinformation
- Privacy
- Lack of context
- Lack of creativity
- Lack of Explainability

Misuse Risks

- Misinformation
- Cyberbullying
- Phishing
- Automated Generation
- Fraud
- Malware

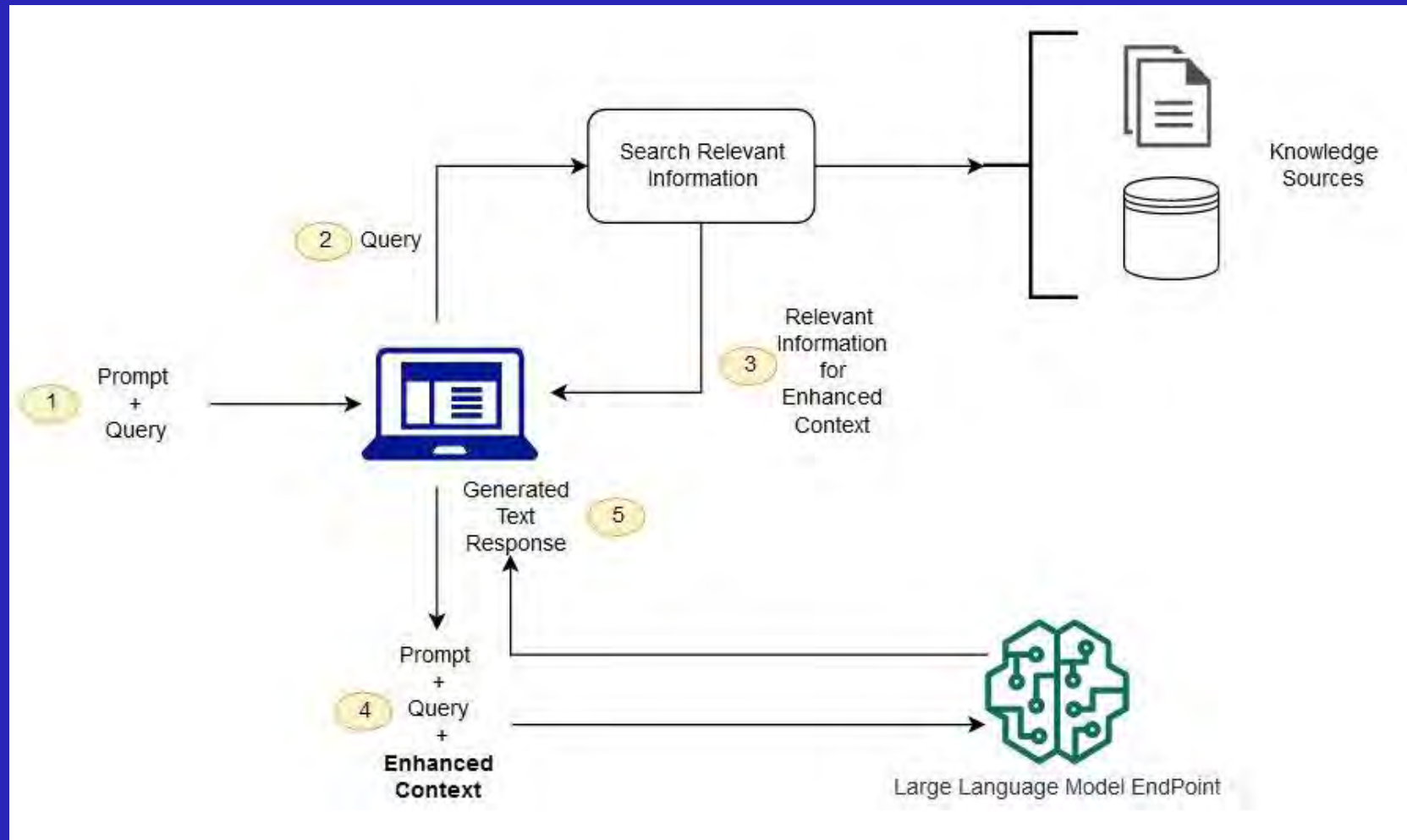
Usage Risks

- Intellectual Property
- Hallucinations
- Copyright

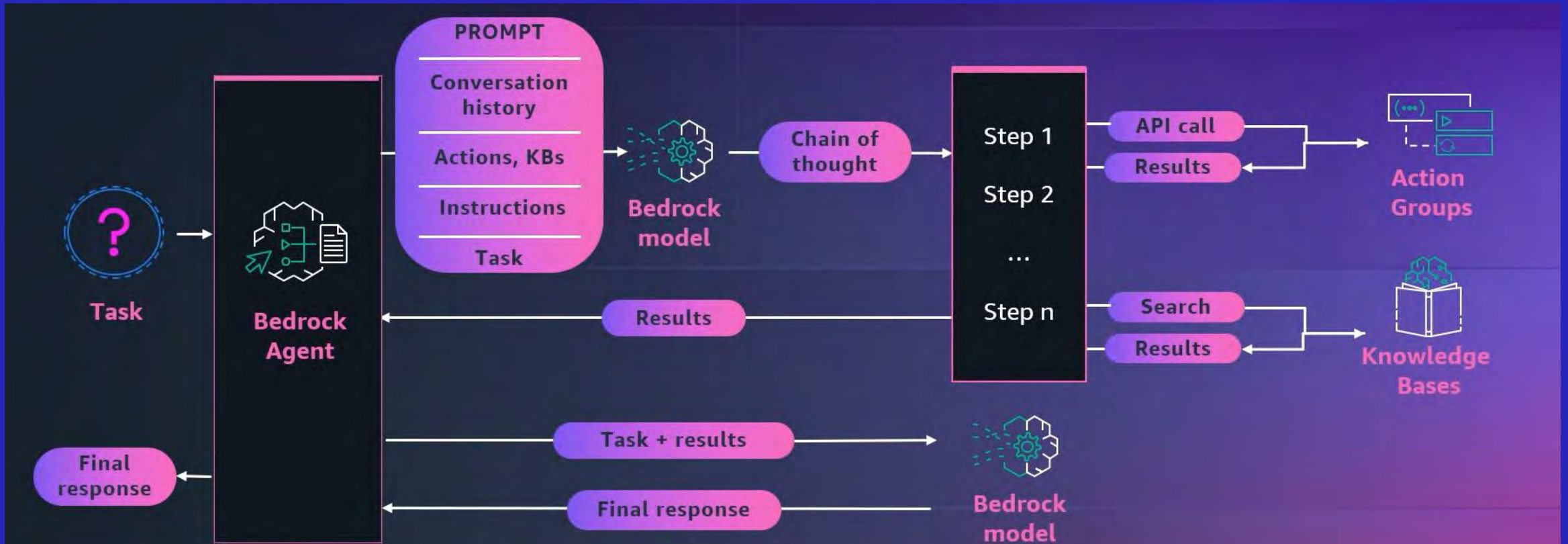
Addressing Model Challenges: Finding Effective Solutions

- RAG / Knowledge bases
- LLM Agents
- Prompt Engineering

Retrieval-Augmented Generation (RAG) / Knowledge bases



LLM Agents



Prompt Engineering Best practices

- **Clear Objectives:** Define clear and specific objectives for each prompt to guide model generation effectively.
- **Relevant Context:** Provide relevant context to the model to ensure accurate and meaningful responses.
- **Continuous Evaluation:** Regularly evaluate prompt effectiveness and adjust as needed to improve model performance.
- **Iterative Refinement:** Continuously refine prompts based on model outputs and user feedback for ongoing improvement.
- **Ethical Considerations:** Consider ethical implications when crafting prompts to promote responsible AI usage.

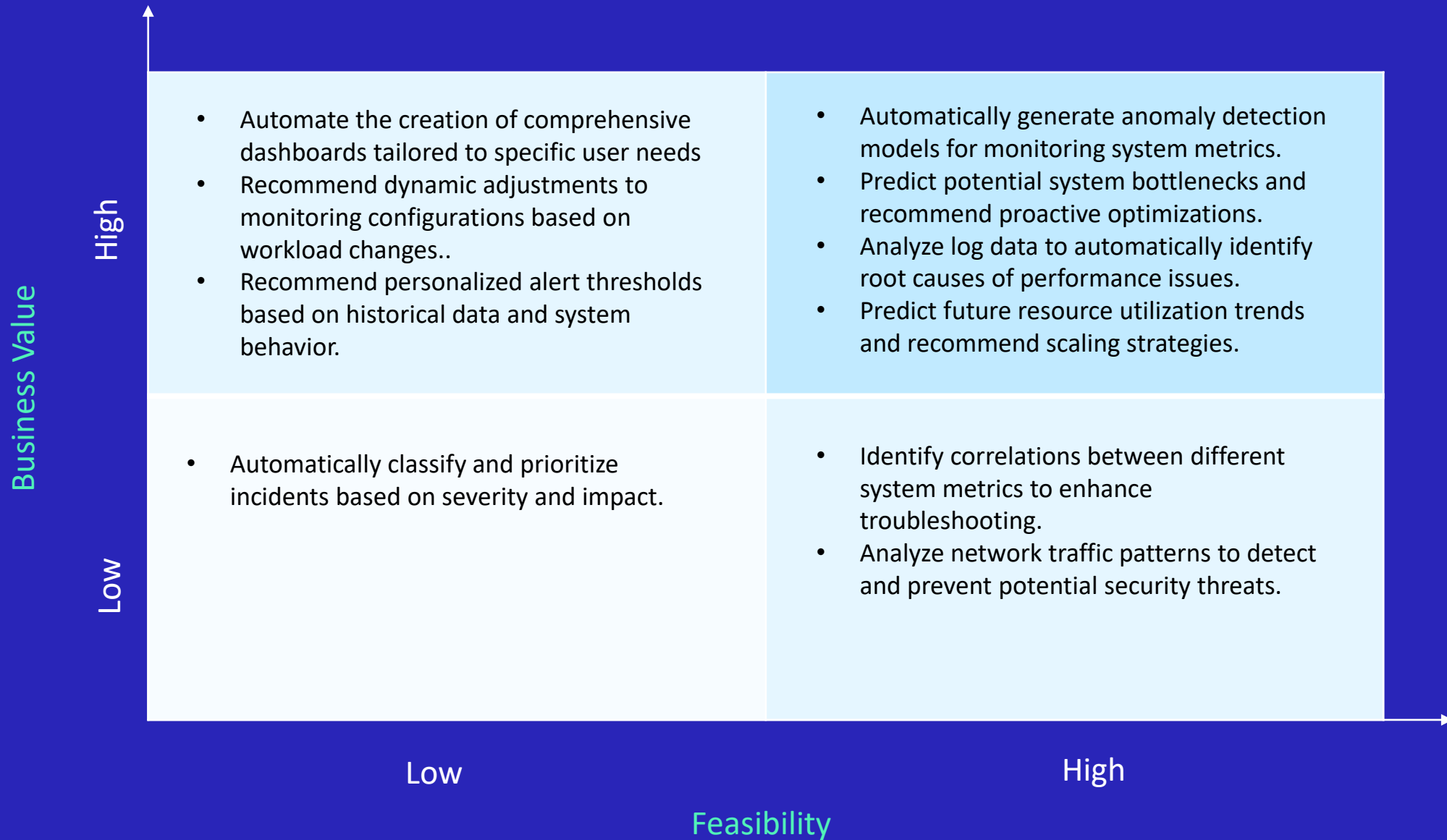
Prompt Engineering Properties

Properties	Details
Temperature	<ul style="list-style-type: none">Controls randomness in model output. Higher temperatures yield more diverse responses; lower temperatures, more focused.
Top-p Sampling	<ul style="list-style-type: none">Controls output diversity by considering only most probable tokens.
Top-k Sampling	<ul style="list-style-type: none">Considers only k most probable tokens for generating next token.
Max Token Length	<ul style="list-style-type: none">Sets maximum length of generated text.
Stop Tokens	<ul style="list-style-type: none">Signals model to stop generating text when encountered.
Repetition Penalty	<ul style="list-style-type: none">Penalizes model for repeating text, encouraging diversity.
Presence Penalty	<ul style="list-style-type: none">Penalizes model for generating already generated tokens.
Batch Size	<ul style="list-style-type: none">Determines number of input sequences processed simultaneously.
Inference Latency	<ul style="list-style-type: none">Time taken for model to generate output given input.
Model Accuracy & Metrics	<ul style="list-style-type: none">Task-specific metrics like accuracy, perplexity, or BLEU score.

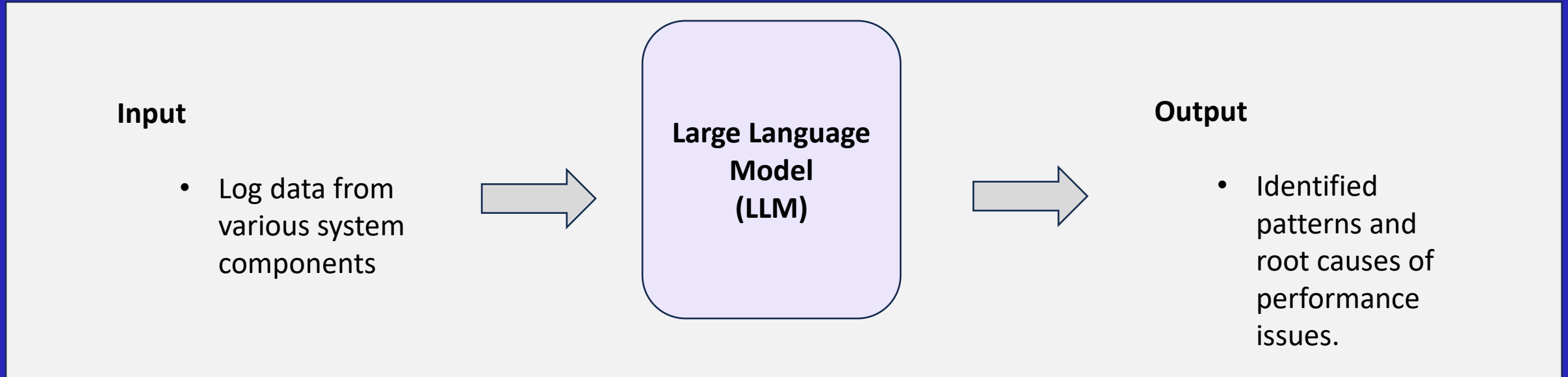
SRE 2.0

- GenAI in Observability
- GenAI in SLI, SLO, and Error Budgets
- GenAI in System Architecture and Recovery Objectives
- GenAI in Release & Incident Engineering
- GenAI in Automation
- GenAI in Resilience Engineering
- GenAI in Blameless Postmortems

GenAI in Observability



Use Case - Analyze log data to automatically identify root causes of performance issues.



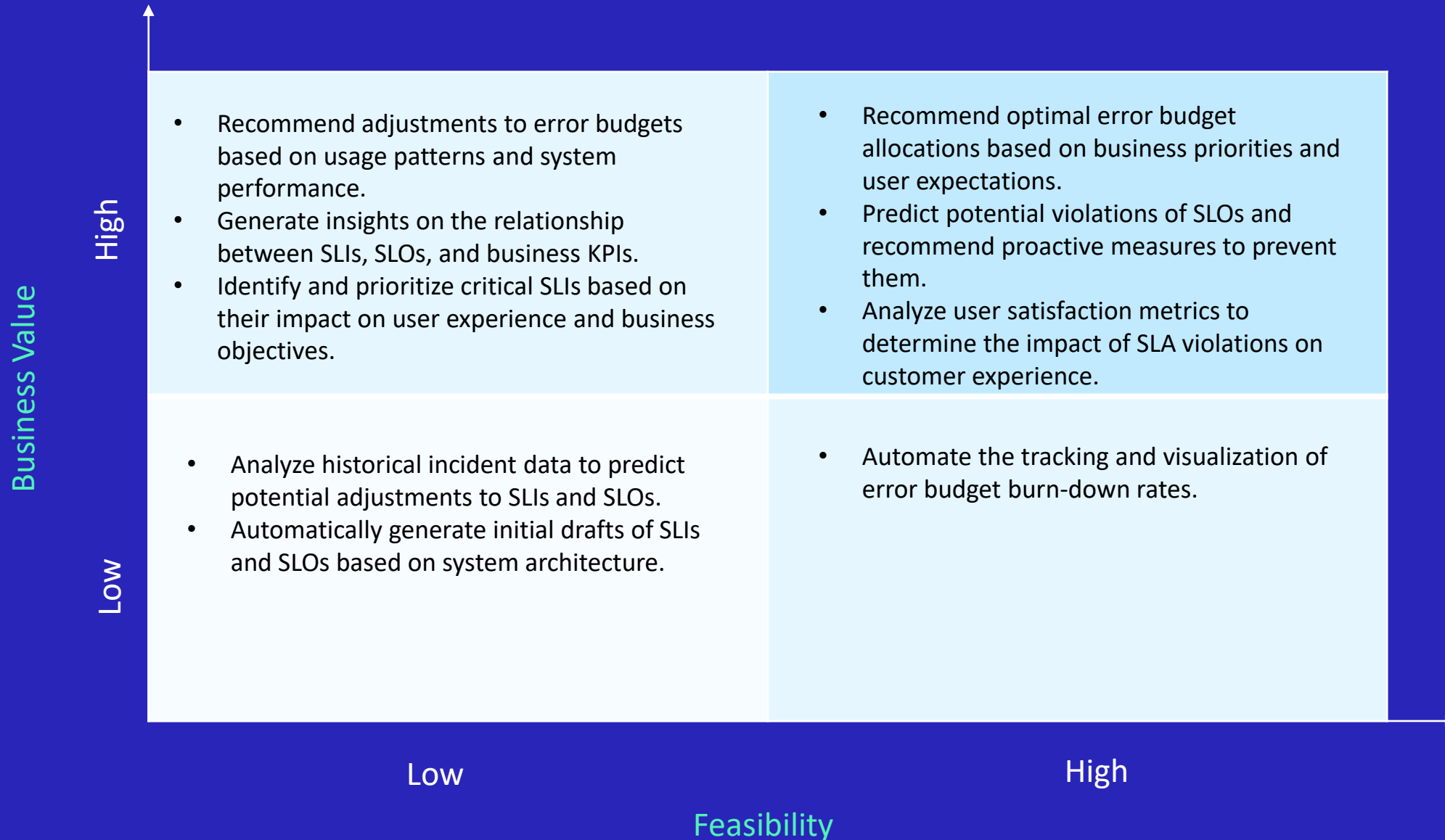
Feedback Loop:

- Human validation of results.
- Feedback from incident resolution.
- Continuous learning.

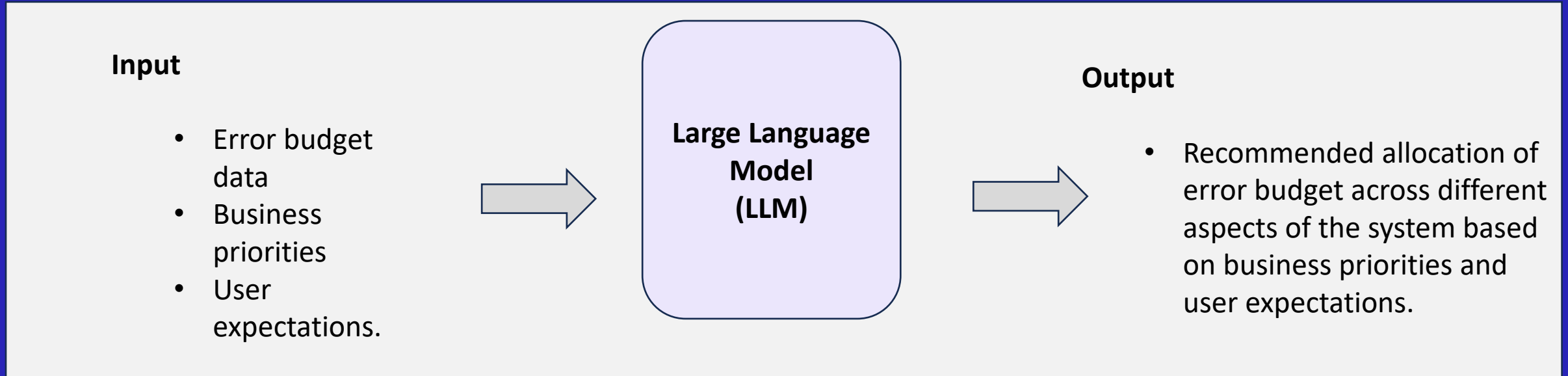
Ways to improve output:

- Enhancing algorithms.
- Providing domain-specific data.
- Integrating feedback effectively.
- Promoting collaborative human-AI interaction.

GenAI in SLI, SLO, and Error Budgets



Use Case - Recommend optimal error budget allocations based on business priorities and user expectations.



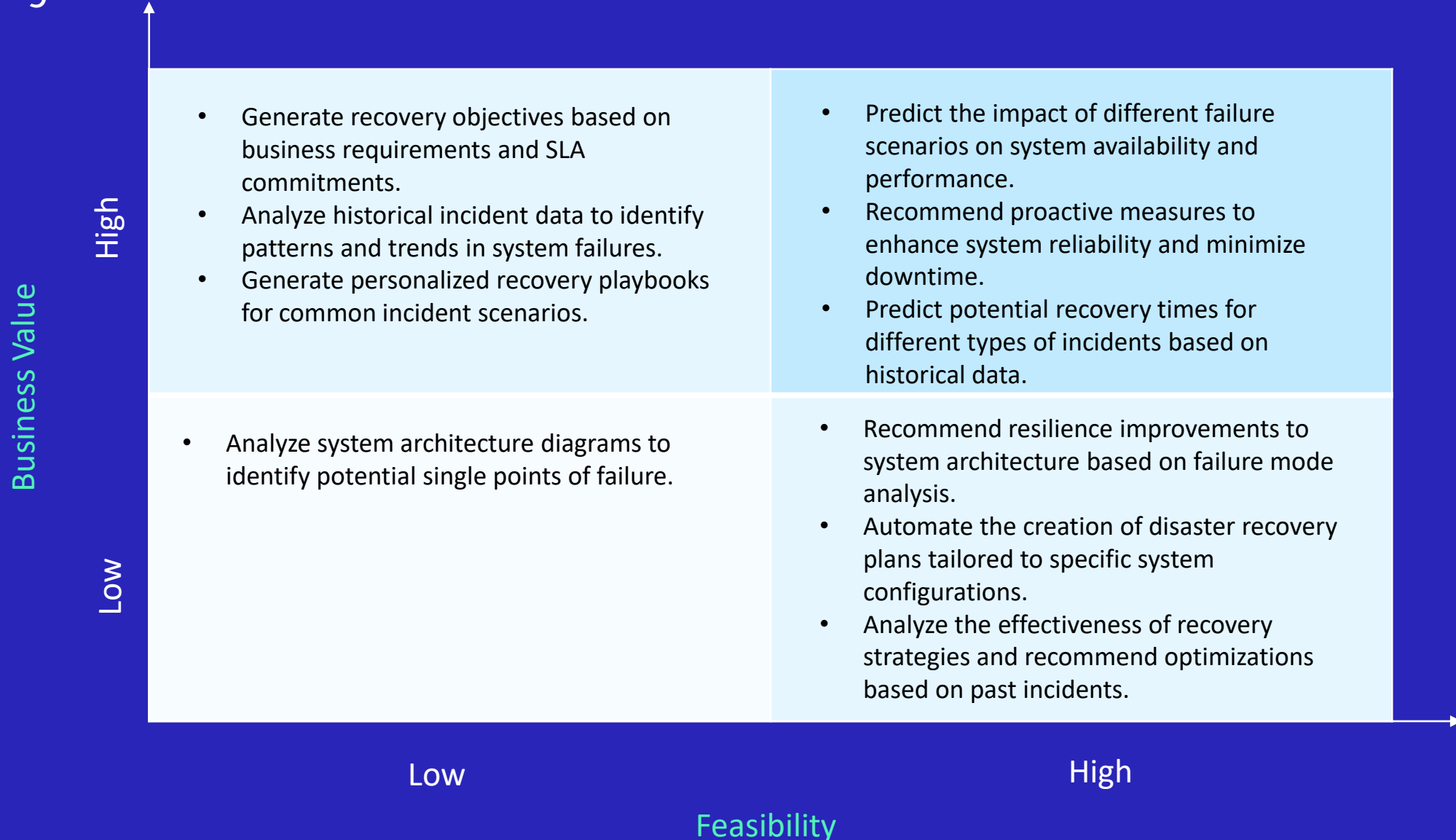
Feedback Loop:

- Evaluation of recommendations by stakeholders.
- Adjustments based on feedback and changing priorities.
- Continuous learning.

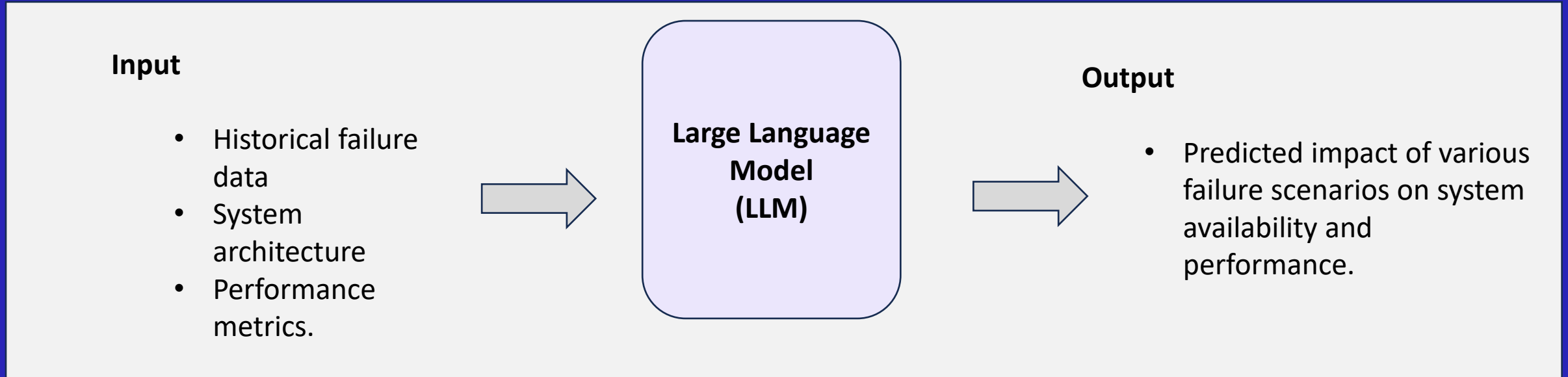
Ways to improve output:

- Incorporating stakeholder feedback into the recommendation process.
- Updating models based on changing business priorities.
- Continuous refinement based on user expectations.

GenAI in System Architecture and Recovery Objectives



Use Case - Predict the impact of different failure scenarios on system availability and performance.



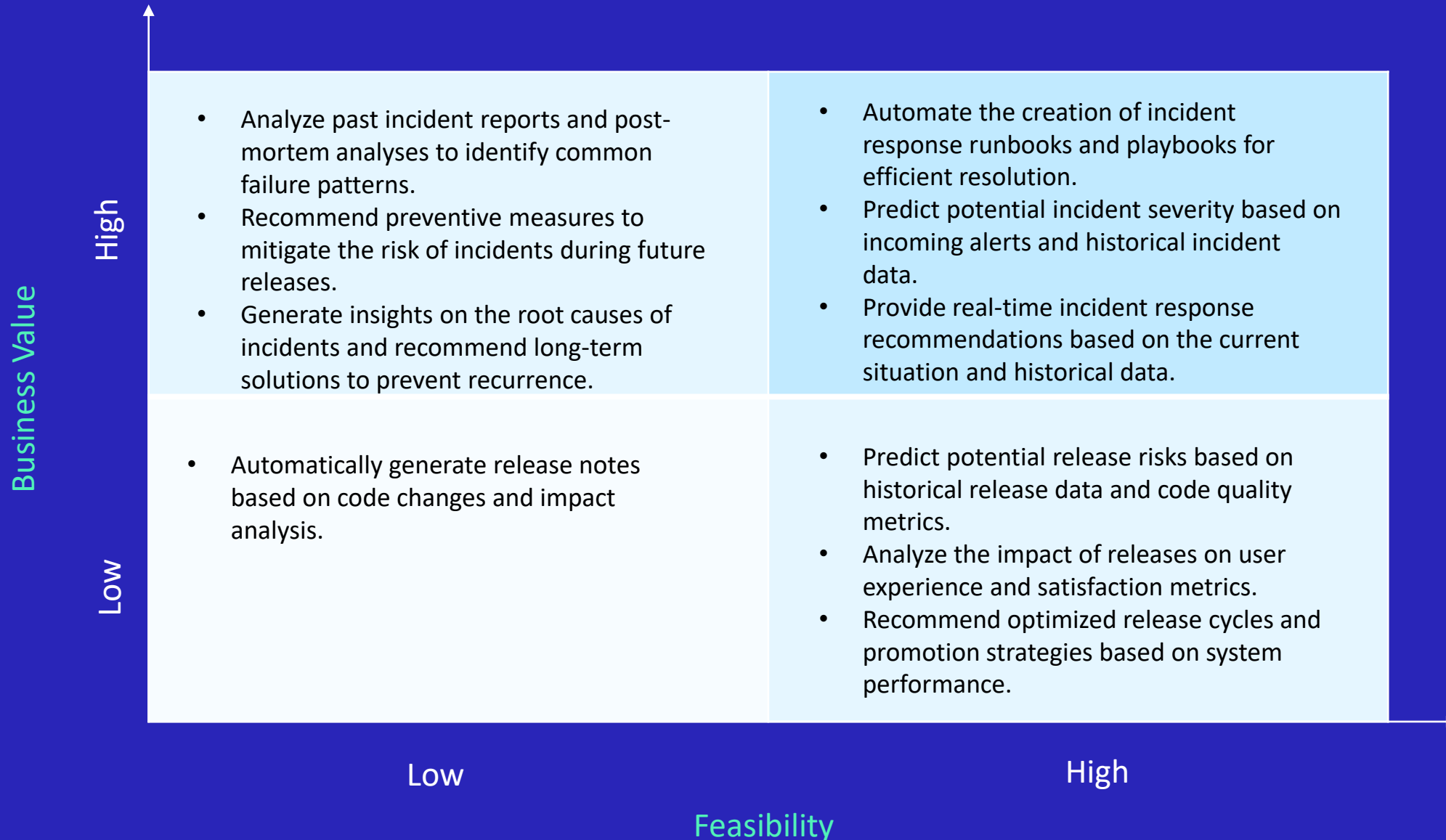
Feedback Loop:

- Validation of predictions against real-world incidents.
- Updating models based on observed impacts and accuracy of predictions.
- Continuous learning from new data.

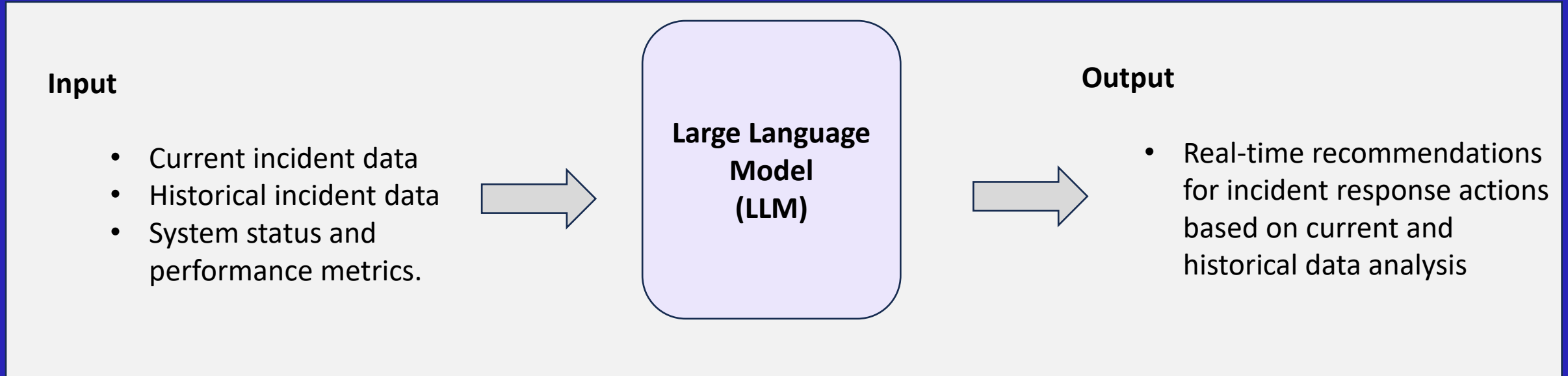
Ways to improve output:

- Incorporating feedback from incident resolution into prediction models.
- Enhancing models with additional factors influencing system availability and performance.
- Continuous refinement based on observed impacts and feedback

GenAI in Release & Incident Engineering



- **Use Case** - Provide real-time incident response recommendations based on the current situation and historical data.



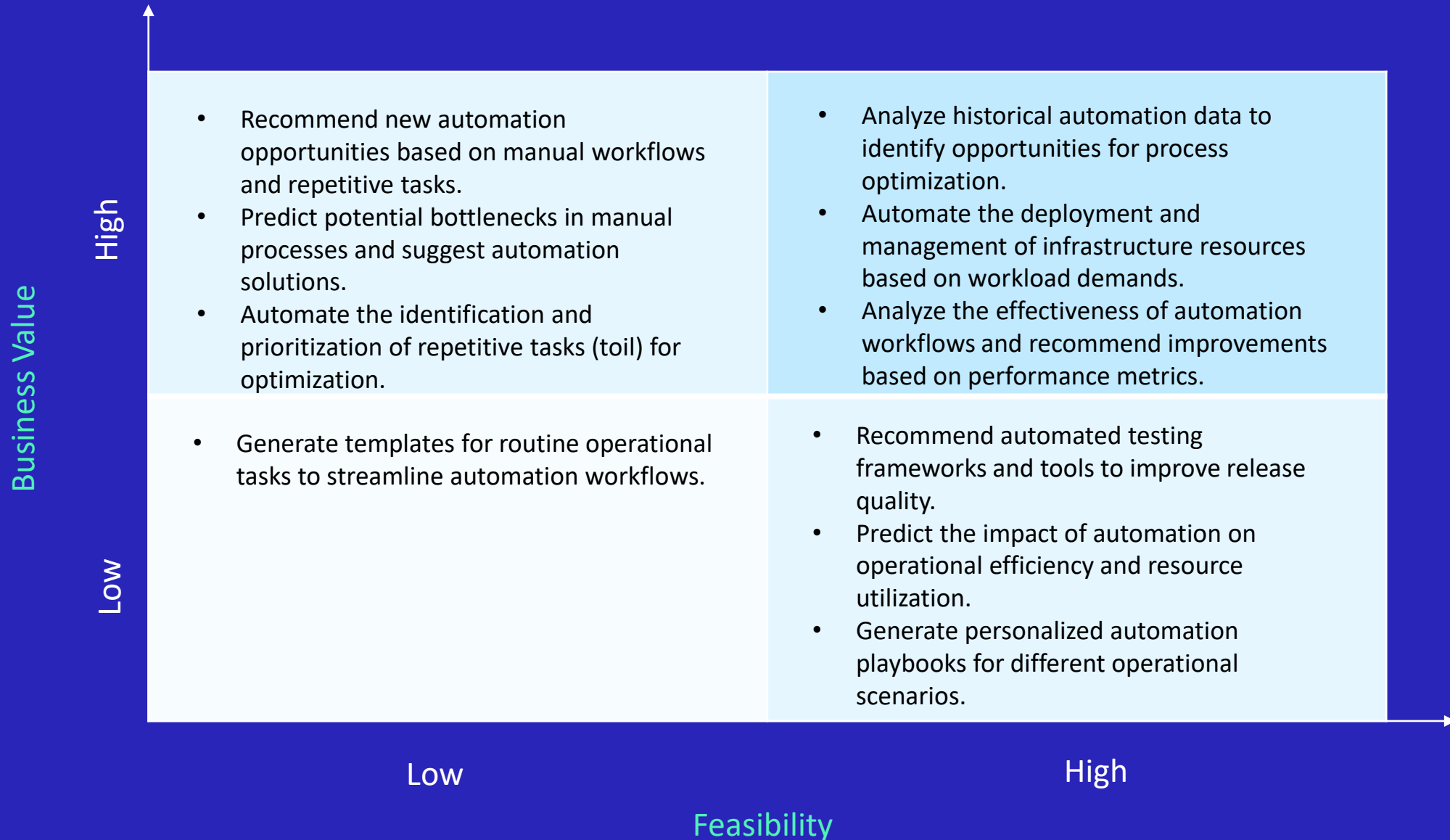
Feedback Loop:

- Evaluation of recommendations' effectiveness in resolving incidents.
- Incorporation of feedback into future recommendations.
- Continuous learning from incident outcomes.

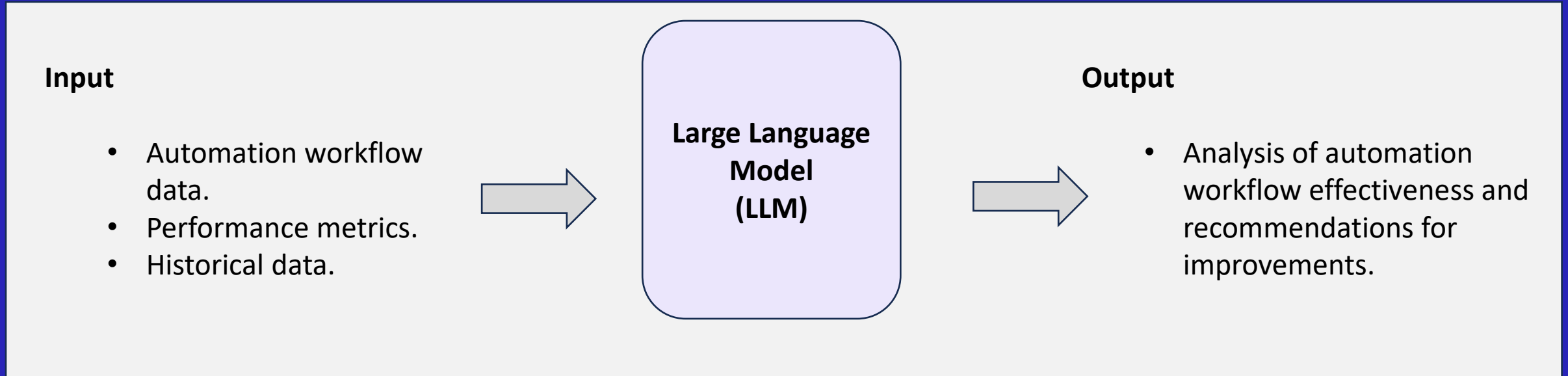
Ways to improve output:

- Integration of real-time feedback from incident resolution into recommendation models.
- Enhancement of models with additional contextual data for more accurate recommendations.
- Continuous refinement based on observed effectiveness and feedback.

GenAI in Automation



- **Use Case** - Analyze the effectiveness of automation workflows and recommend improvements based on performance metrics.



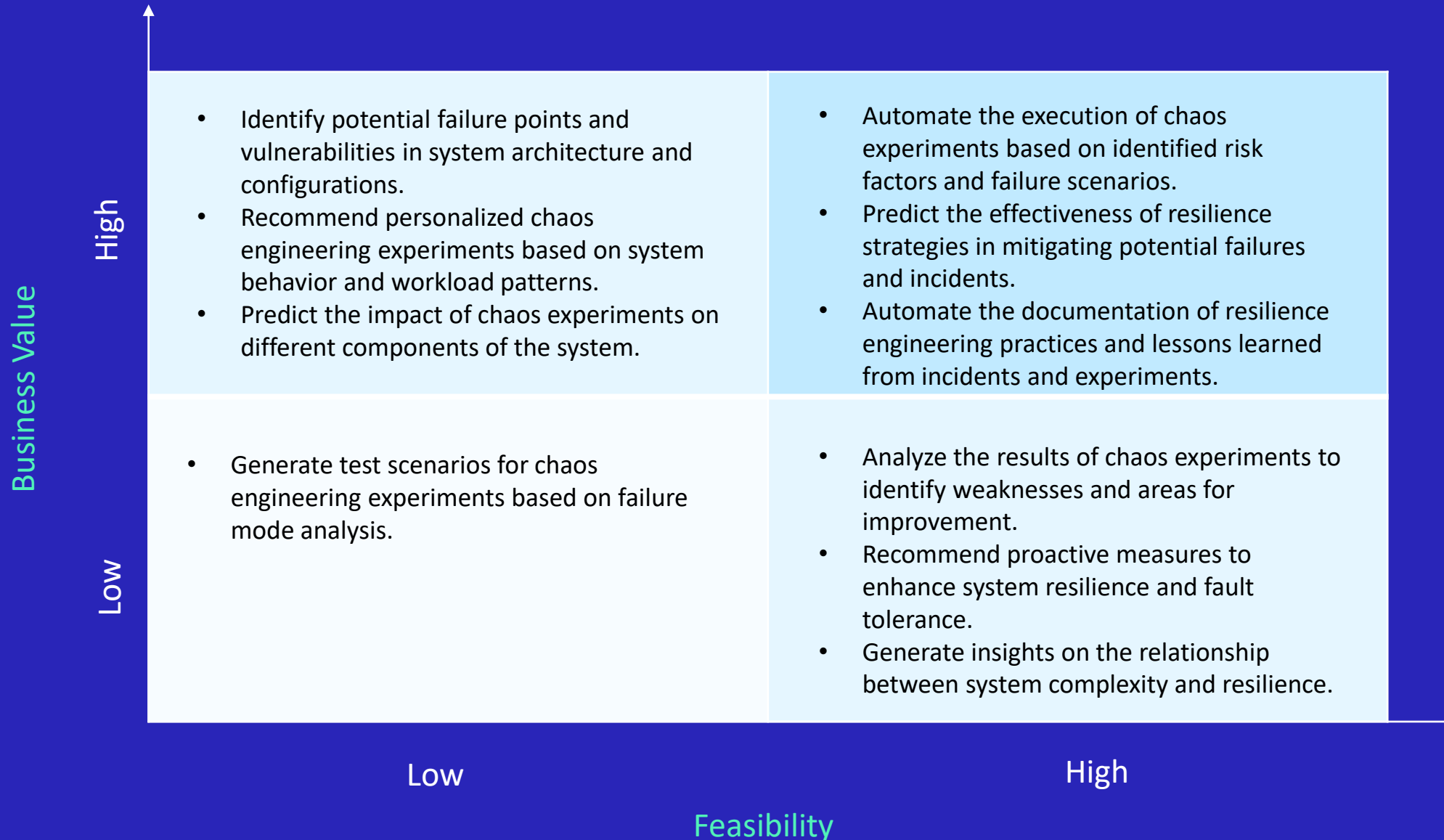
Feedback Loop:

- Evaluation of recommendations' effectiveness in resolving incidents.
- Incorporation of feedback into future recommendations.
- Continuous learning from incident outcomes.

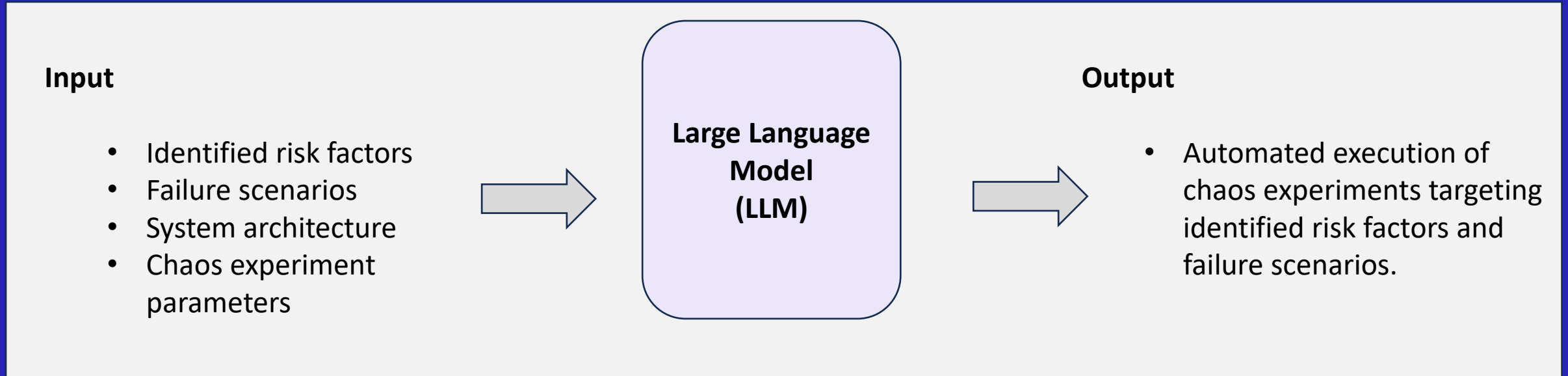
Ways to improve output:

- Evaluation of recommended improvements' impact on workflow performance.
- Incorporation of feedback into future analysis and recommendations.
- Continuous learning from performance outcomes.

GenAI in Resilience Engineering



- **Use Case** - Automate the execution of chaos experiments based on identified risk factors and failure scenarios.



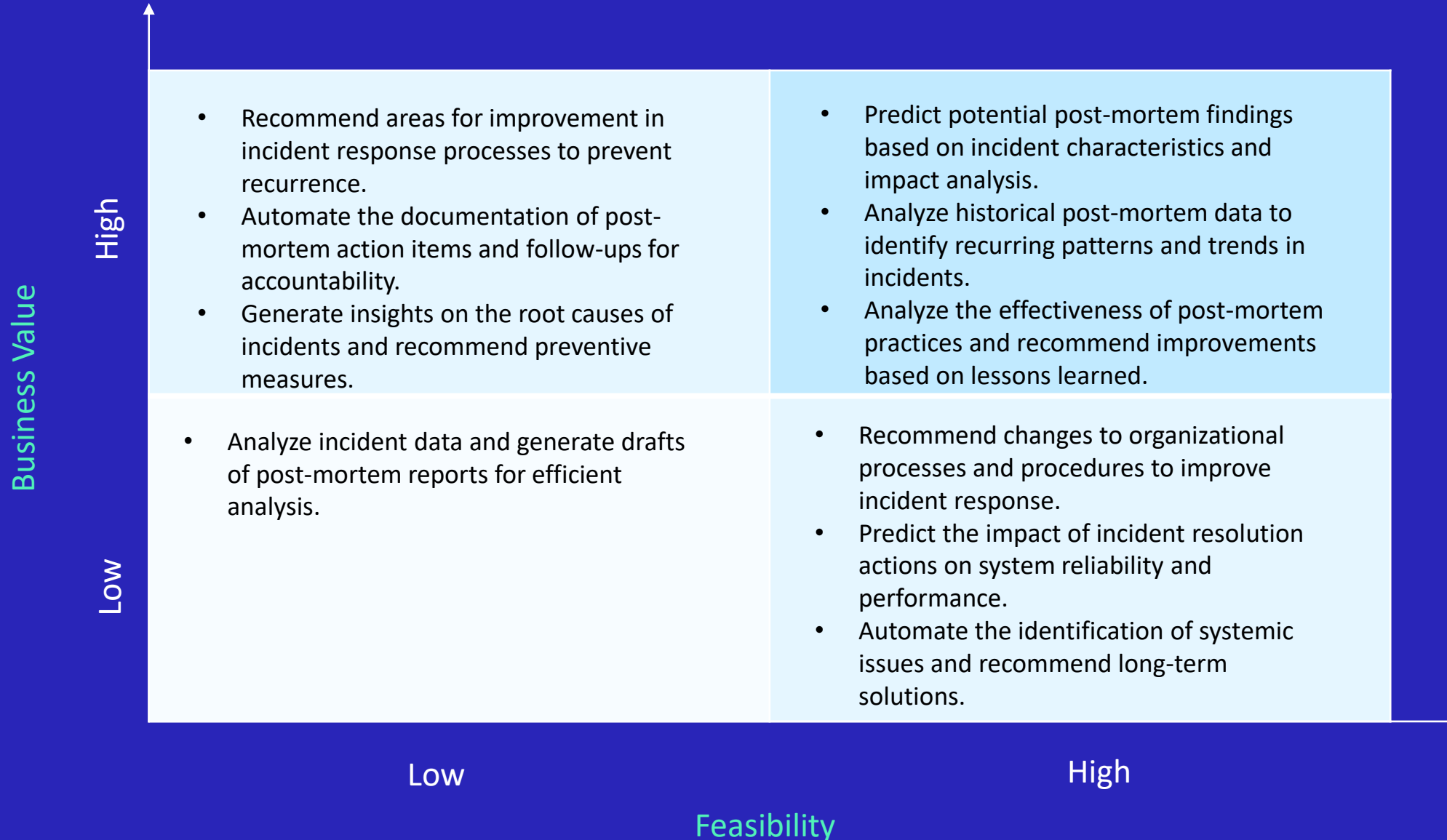
Feedback Loop:

- Evaluation of the impact of chaos experiments on system behavior and stability.
- Incorporation of feedback into future chaos experiment automation.
- Continuous learning from experiment outcomes.

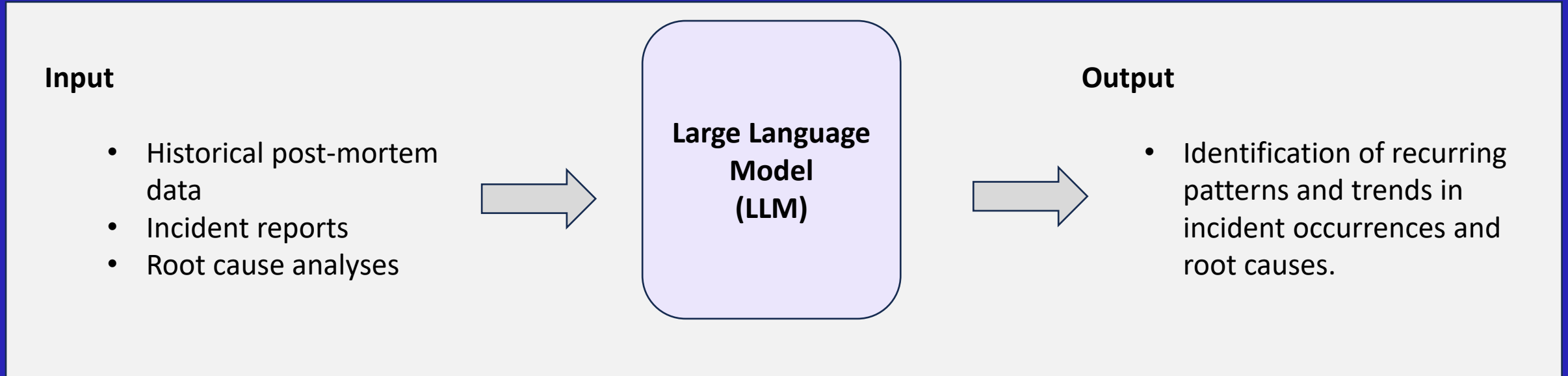
Ways to improve output:

- Integration of real-time feedback from chaos experiments into automation algorithms.
- Enhancement of automation models with additional contextual data for more accurate experiment selection.
- Continuous refinement based on observed outcomes and feedback.

GenAI in Blameless Postmortems



- **Use Case** - Analyze historical post-mortem data to identify recurring patterns and trends in incidents.



Feedback Loop:

- Validation of identified patterns against real-world incidents.
- Incorporation of feedback into future analysis and pattern recognition.
- Continuous learning from incident data.

Ways to improve output:

- Integration of real-time feedback from incident resolution into pattern recognition algorithms.
- Enhancement of analysis models with additional contextual data for more accurate trend identification.
- Continuous refinement based on observed patterns and feedback.

Measure Progress with Business Outcomes

- Improved Alignment with SLOs and Error Budgets
- Increased Change Frequency
- Reduced Change Failure Rate
- Reduced Lead Time for Change
- Reduced Mean Time to Detect (MTTD)
- Reduced Mean Time to Repair (MTTR)
- Increased Mean Time Between Failures (MTBF)



Pitfalls to Avoid

- **Ignoring Ethical Considerations:** Neglecting to address potential biases or ethical concerns in generated outputs can lead to unintended consequences.
- **Lack of Validation:** Failing to validate generated outputs against real-world data or expert judgment may result in inaccurate or irrelevant recommendations.
- **Disregarding User Feedback:** Ignoring feedback from end-users or domain experts can lead to mismatches between generated outputs and actual user needs.
- **Static Models:** Avoid treating generative AI models as static solutions; instead, regularly update and refine them to adapt to evolving requirements and environments.



Thank you.