

Day 12 Model Explainability Guide

A Strategic Framework for Business Leaders
Demystifying AI Decision-Making





Feature importance reveals which inputs have the strongest influence on your model's predictions.

How to Interpret

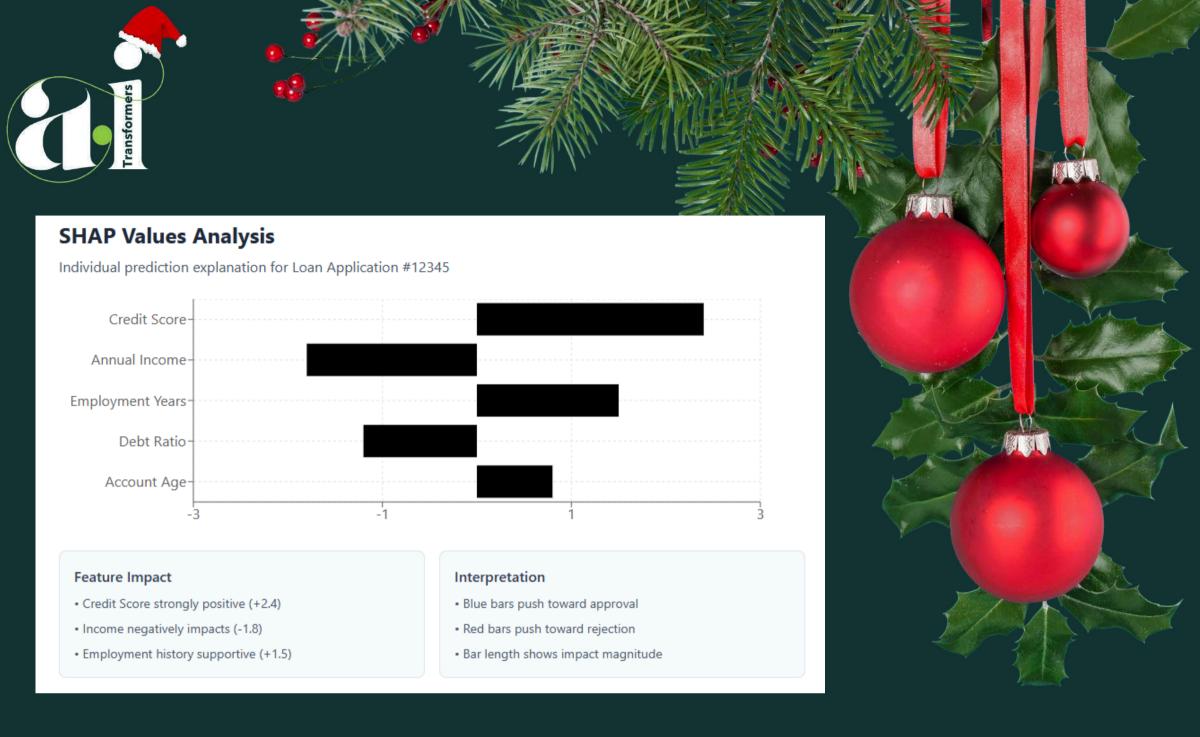
Higher percentages (e.g., 85%) indicate stronger influence on predictions

Look for features that align with domain expertise Question unexpected high-importance features

Red Flags

Single features with extremely high importance (>90%) may indicate oversimplification

Important business variables showing very low importance Unexpected features appearing as highly important



SHAP values explain how each feature contributes to individual predictions, showing both magnitude and direction of impact.

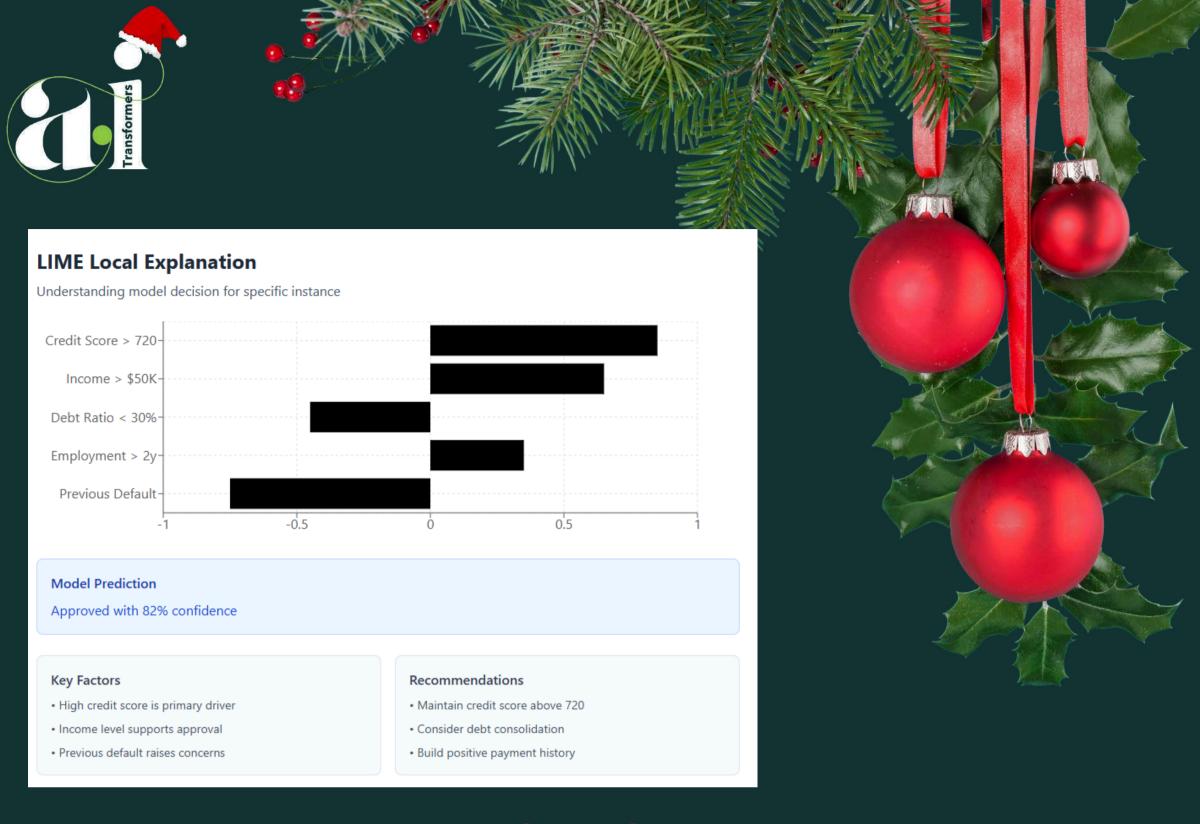
How to Interpret

Positive values: Feature pushed prediction higher Negative values: Feature pushed prediction lower Size of value indicates strength of influence Base value represents average model output



Use for investigating specific decisions
Helpful for regulatory compliance
Valuable for customer-facing explanations





LIME explains individual predictions by creating a simpler, interpretable model around a specific prediction.

How to Interpret

Focus on top contributing factors

Look for logical connections between features and outcome

Compare explanations across similar cases

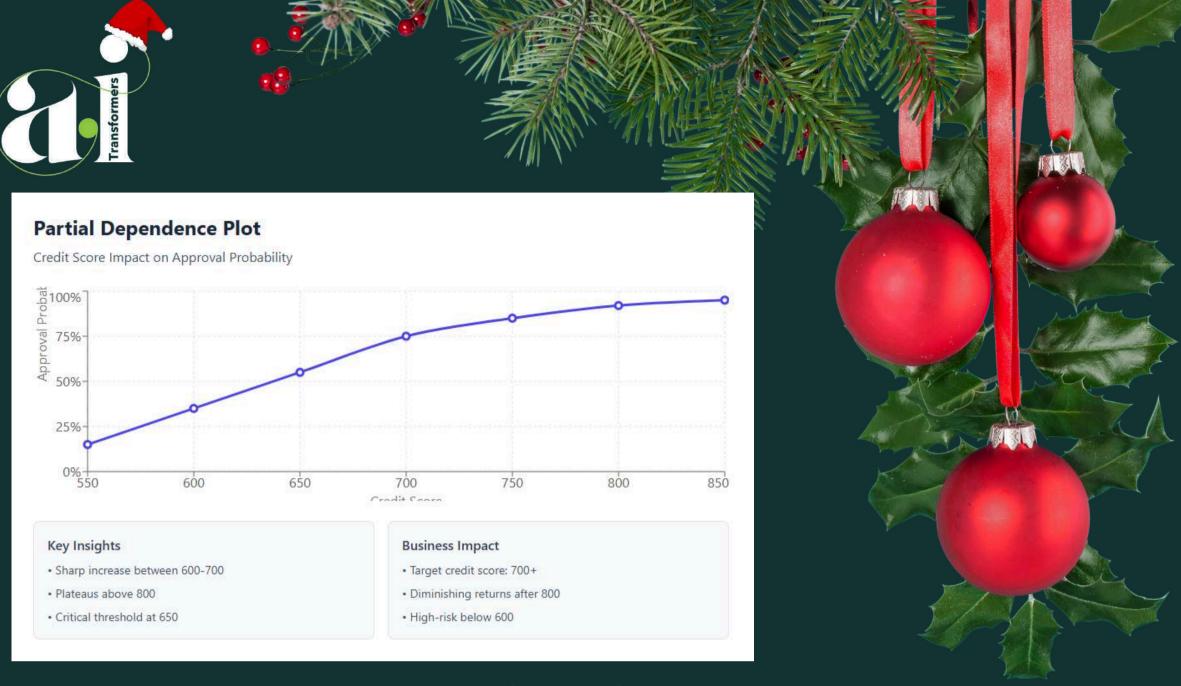
When to Use

Customer service scenarios

Decision appeals

Quality assurance reviews





PDPs demonstrate how changes in one feature affect predictions while holding other features constant.

How to Interpret

X-axis: Feature values

Y-axis: Average predicted outcome

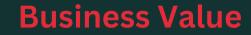
Slope indicates relationship strength

Look for:

Linear relationships

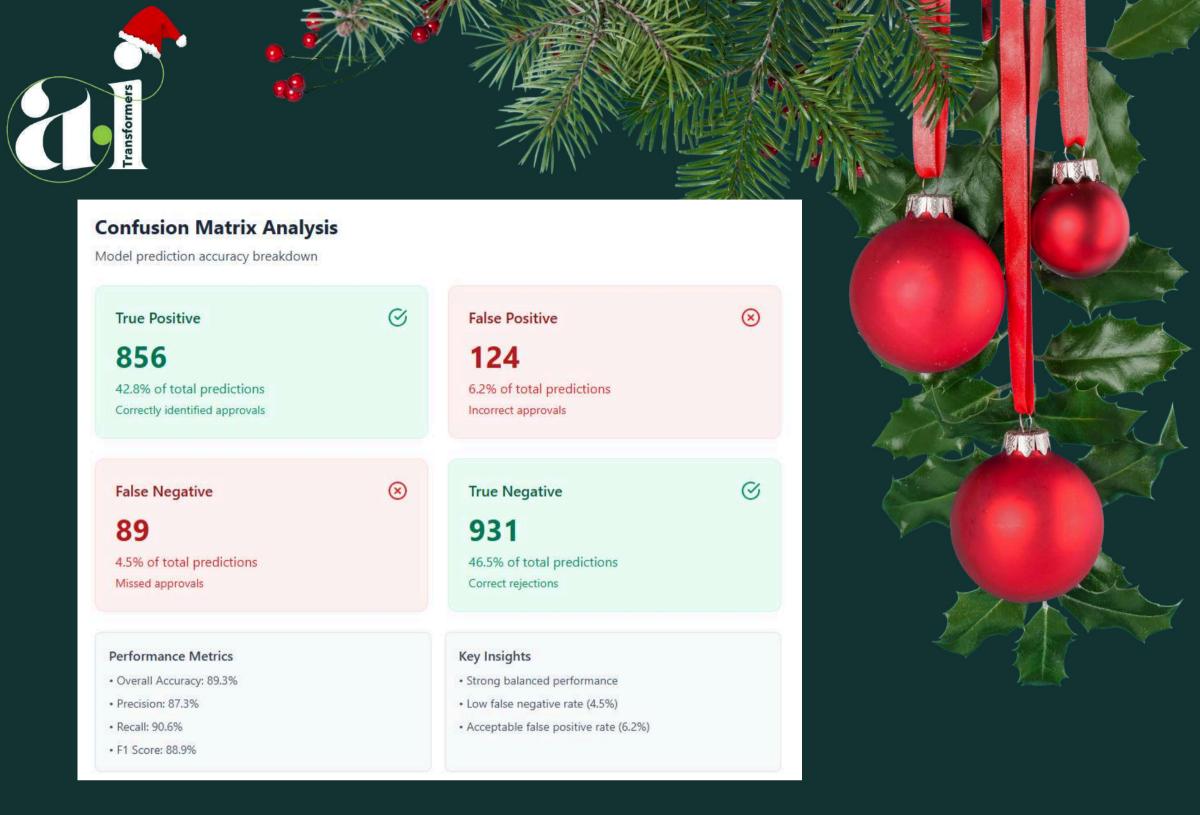
Threshold points

Diminishing returns



Product pricing optimization Risk threshold determination Resource allocation decisions





A comprehensive view of model prediction accuracy, breaking down correct and incorrect decisions.

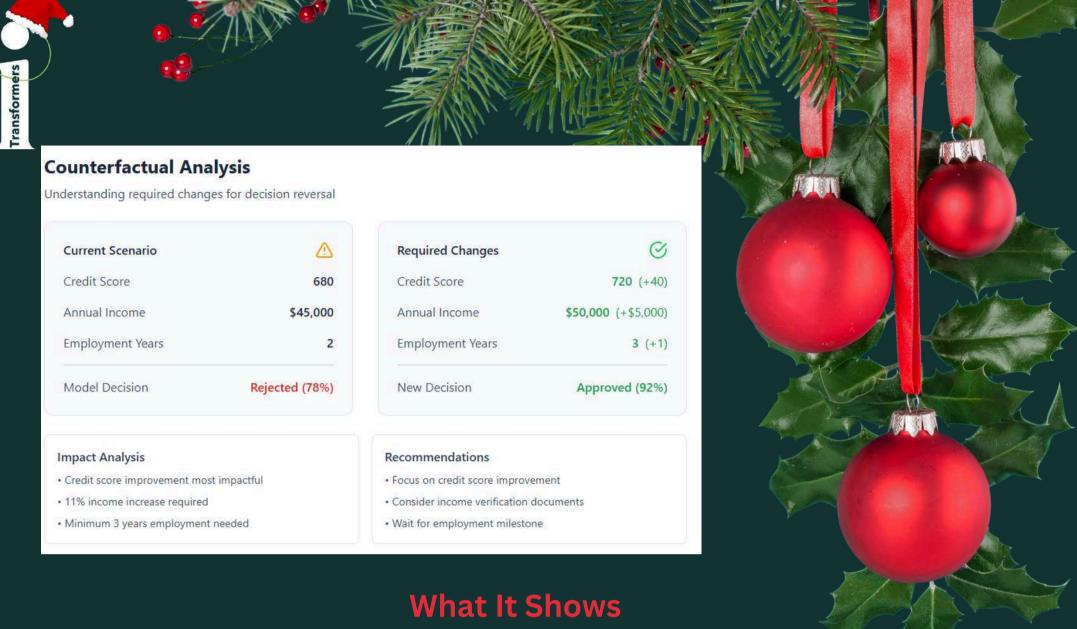
Components

True Positives (TP): Correctly identified positive cases
True Negatives (TN): Correctly identified negative cases
False Positives (FP): Incorrectly identified positive cases
False Negatives (FN): Incorrectly identified negative cases



Business Impact Metrics

Precision: Accuracy of positive predictions
Recall: Ability to find all positive cases
F1-Score: Balance between precision and recall



Counterfactuals demonstrate how input changes would alter model predictions, answering "what-if" questions about model decisions.

How to Interpret

Focus on minimal changes needed to alter predictions Compare against business-feasible scenarios Evaluate practical actionability

Business Applications

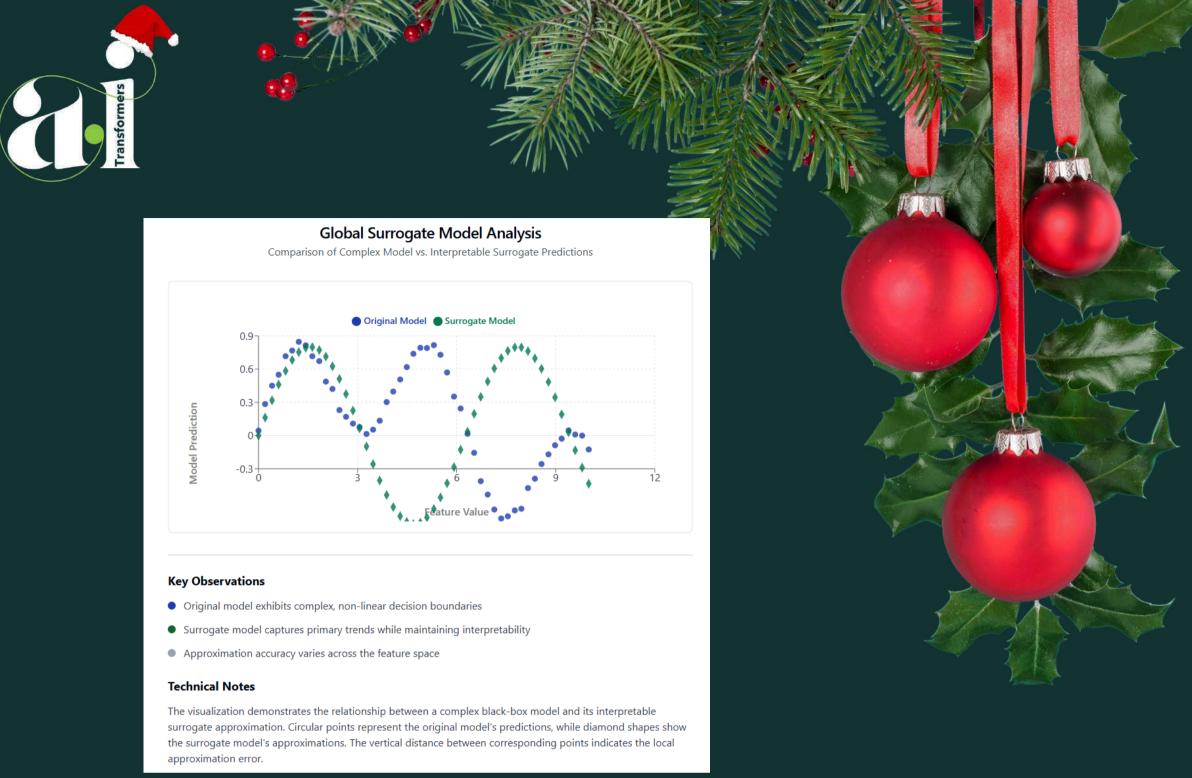
Customer feedback scenarios: "Your loan would be approved if your income were \$5,000 higher"

Process optimization: "Production quality would meet standards with 2°C lower temperature"

Risk management: "Transaction would be flagged as suspicious if amount exceeded \$10,000"

Key Considerations

Feasibility of suggested changes
Cost-benefit of implementing changes
Regulatory compliance implications



Simplified, interpretable models that approximate complex model behavior across all predictions.

Interpretation Framework

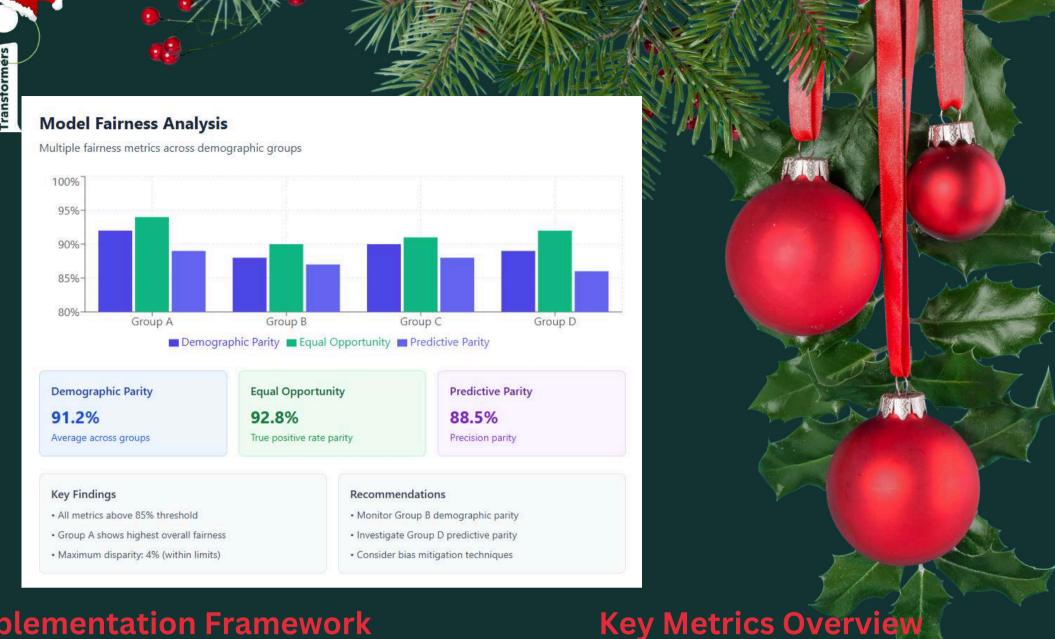
Model similarity score: How well the simple model mirrors the complex one Key decision rules extracted Overall behavior patterns

Business Usage

Executive presentations
Regulatory documentation
Training materials







Implementation Framework

Data Collection Requirements

- Protected attribute identification
- Representative sampling
- Data quality standards

Monitoring Protocol

- Regular fairness audits
- Trend analysis
- Intervention thresholds

Remediation Strategies

- Model retraining triggers
- Bias mitigation techniques
 - Stakeholder communication

Demographic Parity

What: Equal prediction rates groups

When to use: General equity assessment Red flags: Unexplained large disparities

Equal Opportunity

What: Equal true positive rates across

groups

When to use: Performance-based

decisions

Warning signs: Systematic

disadvantages

Predictive Parity

What: Equal precision across groups

When to use: Risk assessment scenarios

Monitor: Group-specific error rates



Integration Best Practices

Holistic Assessment Framework

Metric Combination Strategies

- Complementary metrics selection
- Context-specific weighting
- Composite scoring systems

Governance Structure

- Review frequency
- Responsibility assignment
- Escalation procedures

Documentation Requirements

- Metric baseline establishment
- Change management protocols
- Audit trail maintenance





Regular Reporting

- Executive dashboards
- Operational metrics
- Compliance documentation

Issue Resolution Protocol

- Detection thresholds
- Response procedures
- Communication templates

Continuous Improvement

- Feedback incorporation
- Metric refinement
- Process optimization





- Combine technical and domain expertise
- Verify explanations match business logic
- Challenge counterintuitive results

Documentation Requirements

- Record baseline metrics
- Log significant changes
- Maintain decision rationale

Red Flags and Warning Signs



Model Behavior

- Sudden changes in feature importance
- Contradictory explanations
- Unstable patterns

Business Impact

- Explanations that violate business rules
- Discriminatory patterns
- Counterintuitive relationships



Quarterly Assessment

- Review overall patterns
- Update documentation
- Adjust monitoring thresholds

•

Annual Evaluation

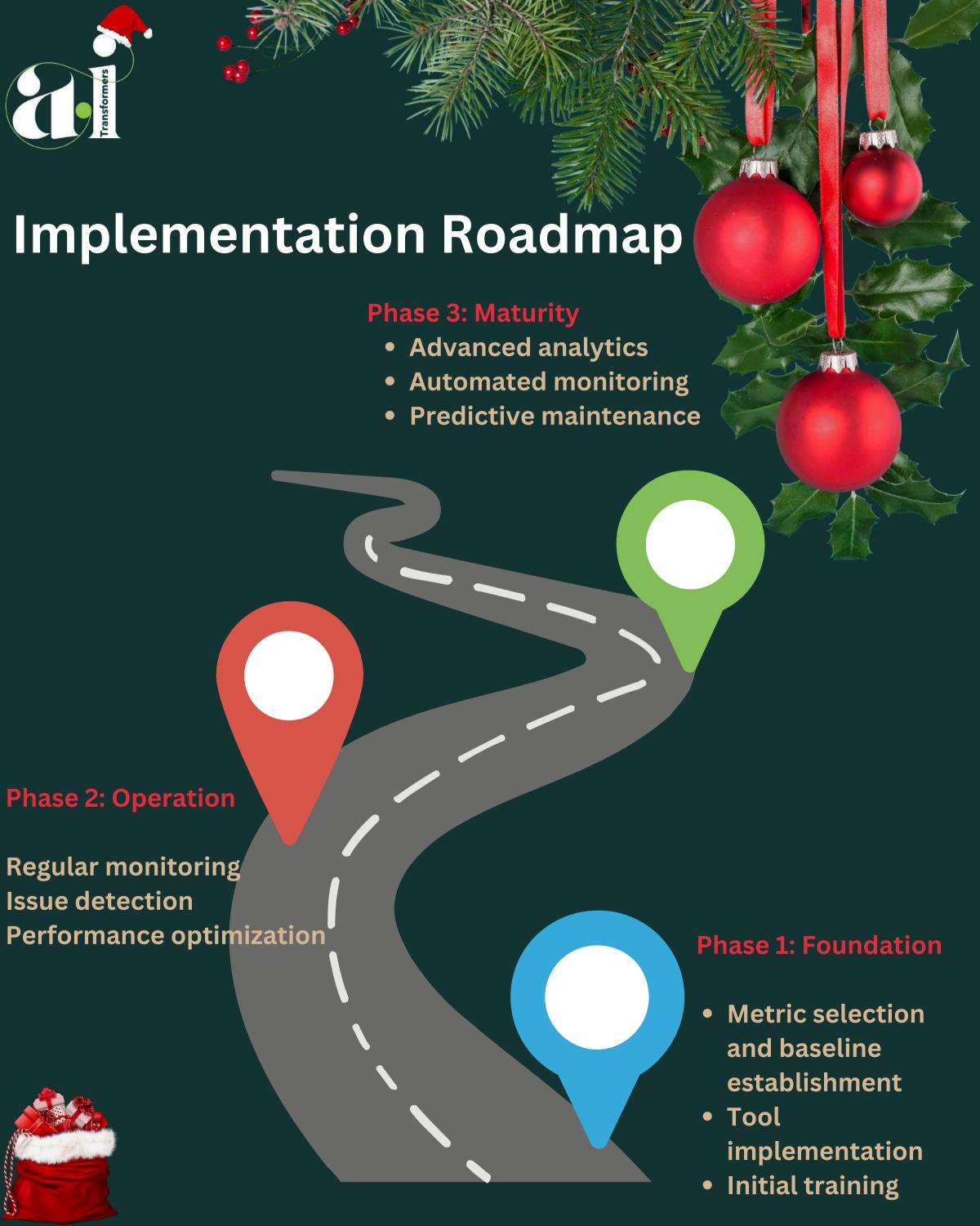
Comprehensive metric review

• Update business alignment

Revise monitoring standards









Transform Your AI Explainability Journey

Partner with Al Transformers to implement robust model interpretation frameworks

Research Impact

92%

Implementation Success

250+

Organizations Analyzed

15+

Industry Sectors

☐ Get in Touch

Discuss your organization's AI explainability needs with our experts.

support@aitransformers.org

⊕ Visit Our Website

Explore our services and research insights.

www.aitransformers.org

☑ Scan QR code for direct access



⇔ Connect on LinkedIn

Follow our latest research and insights.

Al Transformers on LinkedIn

WhatsApp Consultation

Scan to connect with our experts.



Ready to Transform Your AI Explainability?

Partner with AI Transformers to implement robust, interpretable AI systems that drive business value while maintaining transparency and trust.

Contact Our Team

Learn More