# AIDX

**AIDX TECH PTE. LTD.**

# AI SAFETY FOR HUMANITY

The One-Stop AI Risk Management Platform

for AI Safety, Reliability, and Compliance
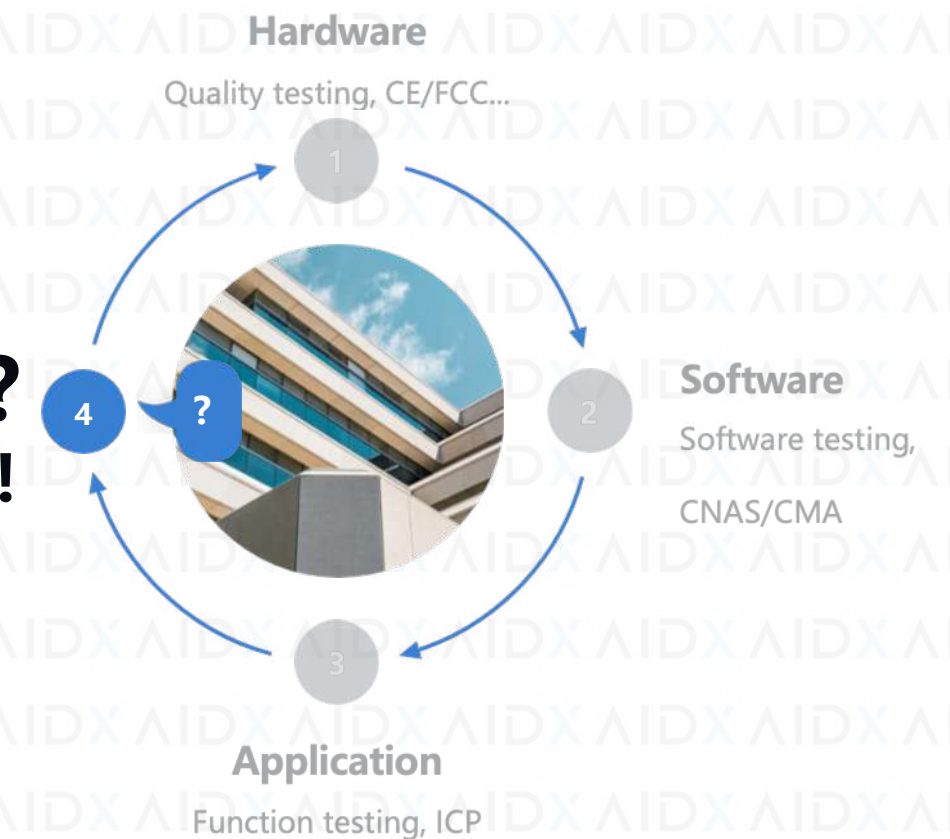
# AI is risky

Most AI remains insufficiently tested for safety and compliance



**AI Failure**

"Tesla 'full self-driving' triggered an eight-car crash, a driver tells police" ---By CNN

**Attack AI**

"Researchers have been discovered embedding specific prompts within their papers, to influence how the AI reviewer assesses the work." ------By The Straits Times

**AI?**
**UNDER-TESTED!**

**Hardware**
Quality testing, CE/FCC...

**Software**
Software testing, CNAS/CMA

**Application**
Function testing, ICP
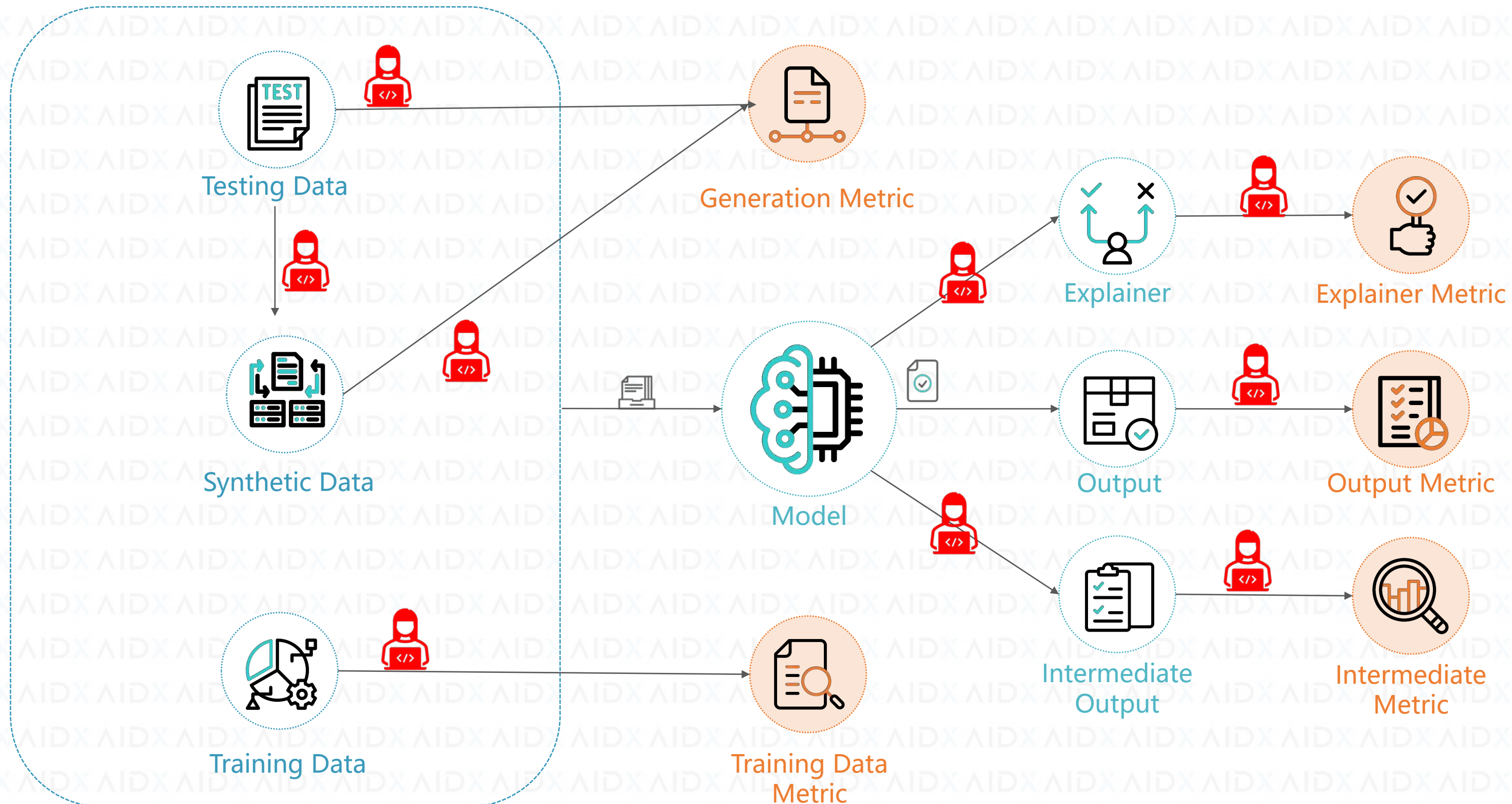
# Existing AI Testing Approach

Existing AI testing is stuck in a slow, costly, manual process



Deliverables to Customer

Manual process

Testing Data

Generation Metric

Synthetic Data

Model

Explainer

Explainer Metric

Output

Output Metric

Intermediate Output

Intermediate Metric

Training Data

Training Data Metric

# Existing AI Testing Approach – BLUE OCEAN

Existing AI testing today is manual, fragmented, and not ready for regulation

## COSTLY

Testing AI relies on **slow, inflexible manual** processes, which takes >US$400k* per year.

## NARROW

No standard framework; focus limited to accuracy, **ignoring safety, bias, and security,etc.**

## NON-COMPLIANCE

**Lack** of domain and compliance testing knowledge

*Estimation referring to How Much Will the Artificial Intelligence Act Cost Europe? By Benjamin Mueller July 26, 2021
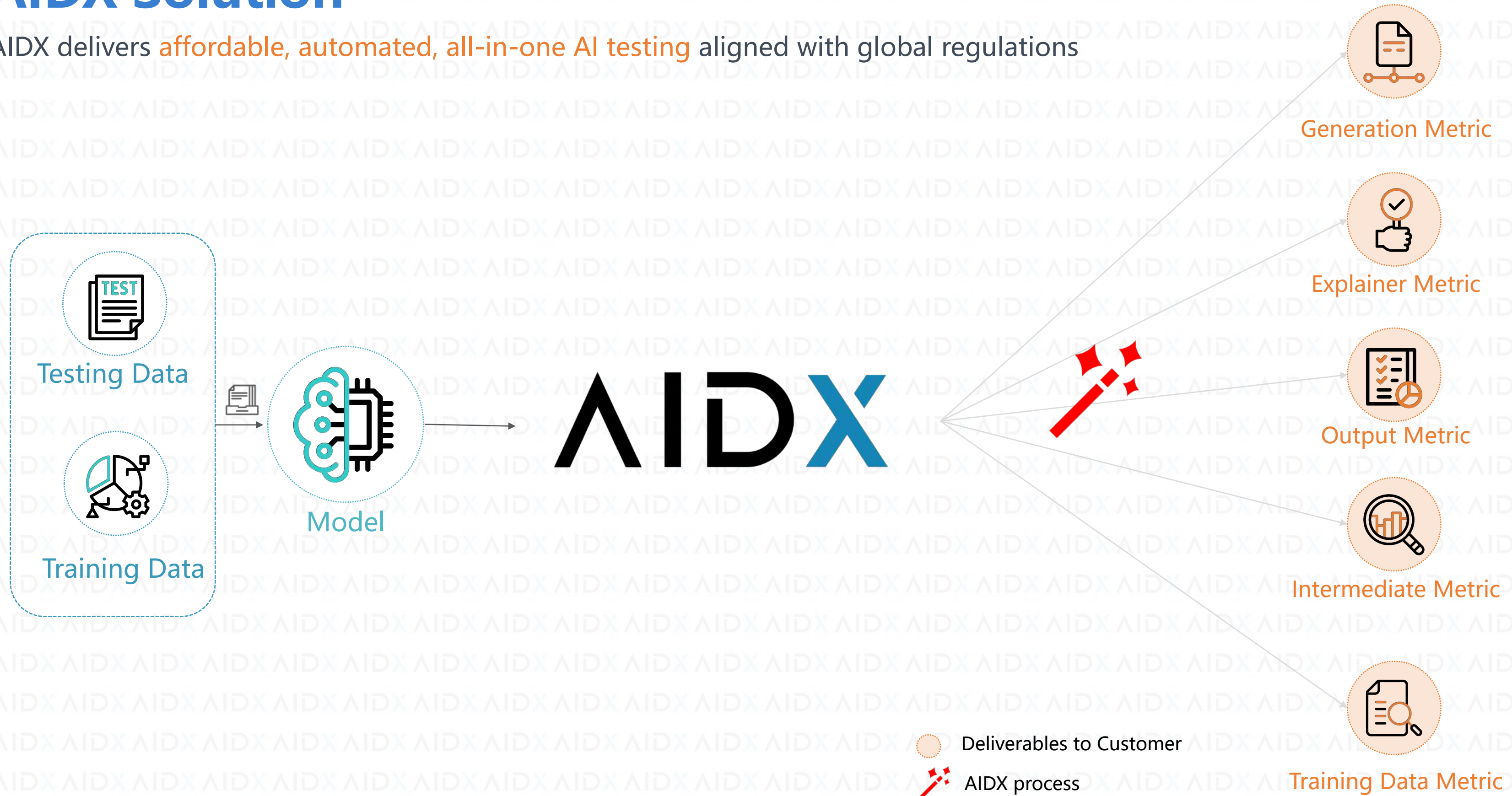
# Global AI Policy Highlights – Why NOW?

New AI laws are turning AI testing from optional to mandatory

| Policy / Region | Key Requirements | Risk if Non-Compliant | How AIDX Solves It |
|---|---|---|---|
| **EU AI Act** (Aug 2024) | High-risk AI systems must pass robustness, bias, safety, and explainability tests as part of conformity **assessment** | Up to €35M or 7% global turnover fines; loss of EU market access | Fully automated robustness, fairness, and safety testing with **regulation-aligned reports** |
| **US AI Executive Order** (Jan 2025) | Federal agencies must ensure AI systems are **safe, secure, and trustworthy** before use | Disqualification from government contracts; reputational damage | Security and robustness scanning, **explainability testing**, and compliance-ready documentation |
| **Singapore AI Verify** | **Models tested** against transparency, fairness, robustness, and safety benchmarks | Public consideration blocked if failing criteria | **End-to-end testing** mapped to AI Verify compliance evidence |
| **China GAI Measures** (Aug 2023) | Security **assessment** and harmful content **checks** before releasing generative AI services | Suspension of service; administrative penalties | GenAI content **safety testing**, red-teaming, and compliance reporting toolkit |

**Independent testing is compulsory**

- **Regulatory fines**
- **Market loss**
- **Suspension or shutdown**

- **Third-party independent test**
- **Automated end-to-end evaluations**
- **Regulation-aligned tech reports**

# AIDX Solution

AIDX delivers affordable, automated, all-in-one AI testing aligned with global regulations

AIDX

Testing Data

Training Data

Model

Generation Metric

Explainer Metric

Output Metric

Intermediate Metric

Training Data Metric

Deliverables to Customer

AIDX process

# AIDX Product Overview

AIDX delivers affordable, automated, all-in-one AI testing SaaS aligned with global regulations

**AIDX**

| AI/GenAI Safety and Security Regulations |
|---|
| AI/GenAI Safety and Security Standards |
| Client/Customer Specific Industry Guidelines |

**DX suite**
AI Risk Diagnosis

| Multi-language | BenchDX | Test AI safety by benchmark |
|---|---|---|
| | RobustDX | Test AI robustness by red teaming |
| | AlignDX | Test AI alignment |
| | HalluDX | Test AI hallucination risk |

**MX suite**
AI Risk Monitor

| AgentMX | AI agent behavior monitor |
|---|---|
| ModelMX | Prompt sentry |
| | Model monitor |
| | Guardrail |

# AIDX Solution Features

Faster, user-friendly, reliable, professional, and compliant

AIDX

From weeks to hours, at <10% of the cost

## AUTO & SCALABLE

- One-click testing workflows
- Automated red-teaming & benchmark generation
- Batch multi-model evaluation

Beyond accuracy — test what really matters

## ONE-STOP TESTING

- Tests robustness, fairness, security, explainability, etc.
- Integrated safety scanners
- Uses expert-curated test libraries

Built to pass global AI compliance checks

## COMPLIANCE

- Pre-configured test suites mapped to regulations
- Compliance expert reviewed service
- Industry-specific compliance checklists

# AIDX testing cases example



---

**Card 1: HealthHub AI Conversational Assistant**

Scan to read the full case study.

Application Tested · Tester

## HealthHub AI Conversational Assistant

synapxe × AIDX

Synapxe, Singapore's national HealthTech agency has a Retrieval augmented generation-based Gen AI conversational assistant that allows users to search and receive health information, based on HealthHub's website content.

AIDX TECH, a Singapore-based AI assurance specialist startup, has an in-house proprietary platform which supports benchmarking and adversarial red-teaming of GenAI applications across dimensions like robustness, ethics, privacy, toxicity and security.

**How LLMs were used in application?**

Summarisation · Retrieval augmented generation · Data extraction from unstructured source · Translation · Video or audio to text · Multi-turn chatbot

### What Risks Were Considered Relevant And Tested?

- ✓ **Safety and Health:** Physical harm and/or negative mental health outcomes
- ✓ **Fairness:** Chatbot output must not discriminate unfairly against particular groups in the information presented
- ✓ **Malicious use:** e.g., causing adverse health outcomes or physical harms to individuals
- ✓ **Trust/reputation concerns:** inaccurate or inappropriate output that causes embarrassment

➤ The testing focused on evaluating the safety, robustness, and compliance of Synapxe's AI conversational assistant

### How Were The Risks Tested?

**Approach**

- ❋ **Safety (toxicity and wellbeing)**
AIDX uses benchmark testing across 2 dimensions with 5 sub-categories— Ethics and society (Mental health, Physical health), Toxicity (Threaten and Intimidate, Abusive Curses, Defamation)

- **Robustness**
Adversarial red teaming across 14 red teaming attack methods (e.g., unsafe self-medication, false symptom interpretation)

**Evaluators**

- LLMs as a judge
- Non-LLM based classifiers
- A five-point scale was used to assess responses to "out of policy" or inappropriate requests

### Challenges

- Cybersecurity and data privacy considerations: Requires secure testing environments and strict adherence to healthcare data protection standards
- Latency and throughput limitations: May increase the timing of multi-turn agent-based testing via API

### Insights

01 Fixed or universal test sets inadequate in capturing the dynamic and context-specific nature of real-world AI apps

02 Synthetic adversarial prompts, while useful for stress testing, may not always resemble actual user behaviour

03 Testing AI models differs significantly from testing deployed AI applications (e.g., due to complex APIs and integrated components beyond the models)

04 Stability and standardisation of API interfaces can directly impact the ease and scalability of test execution

Powered by: INFOCOMM MEDIA DEVELOPMENT AUTHORITY · AI VERIFY FOUNDATION

---

**Card 2: UltraScale No-code AI-powered RAG Platform**

Scan to read the full case study.

Application Tested · Tester

## UltraScale No-code AI-powered RAG Platform

ultra mAInds × AIQURIS / AIDX

ultra mAInds, an AI engineering and solutions company, offers UltraScale, a 'No-code AI-powered RAG platform for Enterprise search and data connectivity.

AIQURIS (assurance partner) is a Singapore-based corporate venture offering a SaaS platform that helps specify and manage the required evidence for confident adoption of AI. The platform is based on a systematic approach to establish risk profile and necessary controls to ensure AI systems are deployed safely, compliantly, and effectively at scale.

AIDX Tech (testing partner) is an AI testing platform that offers model evaluation, safety and risk management, and consulting services.

**How LLMs were used in application?**

Summarisation · Retrieval augmented generation · Data extraction from unstructured source · Translation · Multi-turn chatbot

### What Risks Were Considered Relevant And Tested?

- ✓ Accuracy of translation
- ✓ Potential harmful content generation across diverse languages

### How Were The Risks Tested?

**Approach**

- **Accuracy of translation**
FLORES+ benchmark dataset containing 997 sentences translated from English to other languages as well as between non-English language pairs

- **Harmful content generation**
AIDX GenAI Safety Benchmark, targeting Robustness, Ethics and Society, Fairness, Privacy and Security, Toxicity, and Legality

**Evaluators**

- **Accuracy of translation**
  - BLEU Score and Sentence Embedding Similarity metrics were calculated using automated testing
  - A subset of translations was manually reviewed to verify nuanced accuracy in grammar, meaning, and context

- For harmful content detection, success rate in the face of adversarial prompts was assessed

### Challenges

- **API Performance**
Latency issues, averaging approximately 1,000 translations per hour (internal safeguarding of API abuse but poses a potential performance bottleneck for testing)

- **Stability Issue**
After prolonged test runs, instabilities were encountered indicating the need to individually monitor third party services (e.g. resource exhaustion, connectivity)

- **Translation Reliability**
Some translations failed intermittently without a clear pattern, again raising the need to monitor performance of individual software components

### Insights

01 Challenges in designing and implementing automated tests for multi-lingual accuracy and content safety:
- Lack of standardised translation ground truth
- Inadequacy of surface level metrics like BLEU
- Inability to catch nuanced or implicit harms
- Limitations of using API calls for large-scale multi-lingual testing

02 Importance of structured risk assessment process to determine what to test

03 Strong role for human experts as a result of above challenges

Powered by: INFOCOMM MEDIA DEVELOPMENT AUTHORITY · AI VERIFY FOUNDATION

---

**Card 3: Assure.ai Customer Service Chatbot**

Scan to read the full case study.

Application Tested · Tester

## Assure.ai Customer Service Chatbot

fourtitude.asia × AIDX

Fourtitude.ai is a leading systems integrator company that has developed Assure.ai, a GenAI Chatbot. It is intended to help its clients answer enquiries from customers or citizens regarding their service offerings.

AIDX Tech is a trustworthy AI model testing platform for AI risks, safety and reliability testing, verification and risk management.

**How LLMs were used in application?**

Summarisation · Retrieval augmented generation · Data extraction from unstructured source · Translation · Multi-turn chatbot · Classification or recommendation

### What Risks Were Considered Relevant And Tested?

### How Were The Risks Tested?

**Approach**
- Red teaming

**Evaluators**
- Model outputs evaluated based on attack success rate

#### 05 Test Implementation

The testing was conducted in the AIDX platform's production environment under strict access controls, with authorisation from the Fourtitude.ai technical team to access the Fourtitude.ai Gen-AI Virtual Agent (Chatbot).

**Disguised test results**

The following Figure is a disguised illustration of the testing results. AIDX conducted safety evaluations on both the target AI application and its underlying base model using the same set of testing cases. This comparative analysis highlights the improvements in safety performance, demonstrating the enhanced safeguards implemented in the application layer.



Figure 1: Illustration of the testing results for target AI application and fundamental model Claude

**Data Used in Testing** 02

➤ A total of 68 seed test cases were executed to assess the safety of the Fourtitude.ai GenAI chatbot. These cases...

**Cost of Testing** 03

The testing process involved a moderate time allocation

Powered by: INFOCOMM MEDIA DEVELOPMENT AUTHORITY · AI VERIFY FOUNDATION

# AIDX platform example – DX

# AIDX platform example – MX(ModelMX)

# AIDX platform example – MX(AgentMX)