

# AI Explainability WG

May 2023

# May 2023 updates

- TrustyAI service
  - Improvements to metrics endpoints
    - Calculation data size configurable per-metric
    - Metadata includes representative values
    - More robust REST payload validation
    - Per-metric configurable thresholds
- Python TrustyAI 0.2.12
  - Mostly fixes around the Tyrus visualisations<sup>[1]</sup>
- Community
  - Architectural Decision Records (ADRs) for proposals<sup>[2]</sup>

[1]: <https://pypi.org/project/trustyai/0.2.12/>

[2]: <https://github.com/trustyai-explainability/community/tree/main/adr>

Community

# Community

- Community repository<sup>[1]</sup>
  - Community information
  - Contribution guidelines
  - Code of conduct
- Architectural Decision Records (ADRs)
  - Approved ADRs:
    - ADR-0001: TrustyAI external library integration<sup>[2]</sup>
    - ADR-0002: Metrics and XAI namespaces<sup>[3]</sup>

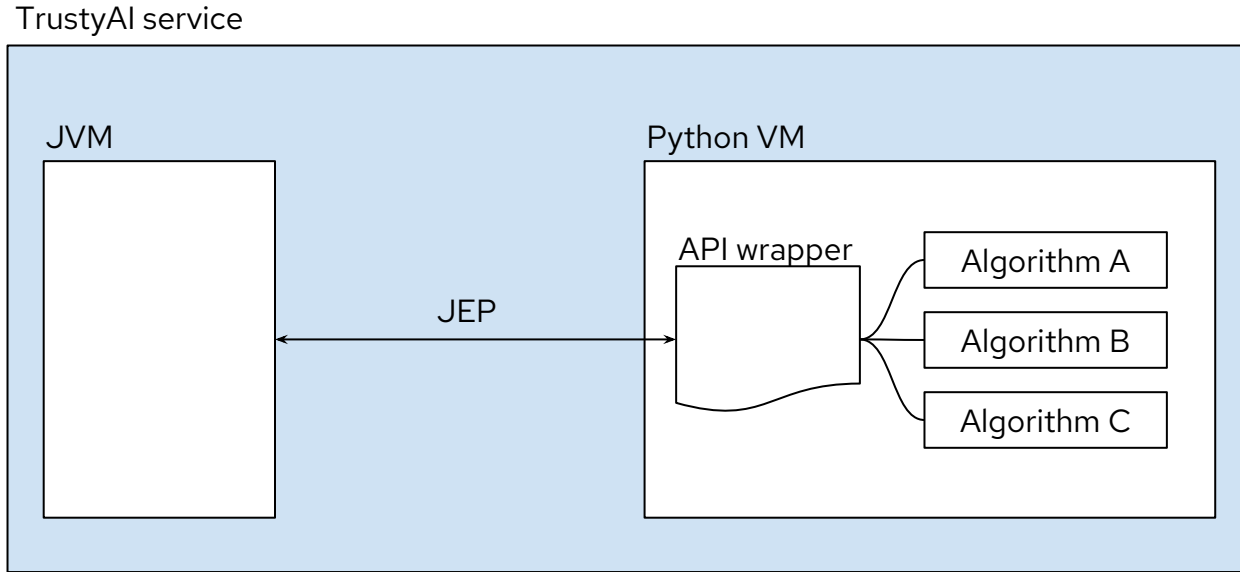
<sup>[1]</sup>: <https://github.com/trustyai-explainability/community>

<sup>[2]</sup>: <https://github.com/trustyai-explainability/community/blob/main/adr/ADR-0001-trustyai-external-library-integration.md>

<sup>[3]</sup>: <https://github.com/trustyai-explainability/community/blob/main/adr/ADR-0002-metrics-and-xai-namespaces.md>

TrustyAI external library integration

# External library support



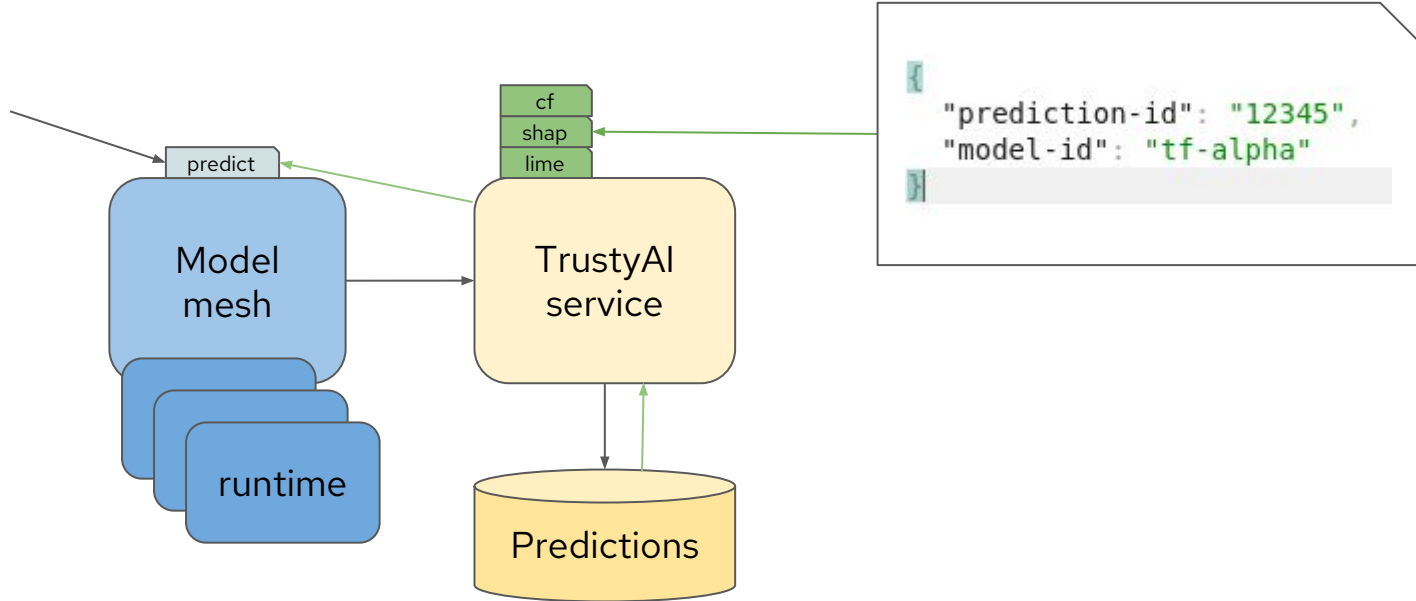
# External library support

- Communication API
  - Ongoing development
  - Determine general requirements for metrics / explainers
    - Representative APIs
    - Data structures for marshalling
    - Performance benchmarks

# TrustyAI Explainers Endpoints



# Explainers Endpoints



# Roadmap

# TrustyAI 2023 roadmap

## July 2023 TrustyAI 0.2.0

- *Explainers*
  - Support for explainers LIME, SHAP, CF at service level
- *Metrics*
  - Flexible scheduling/batching
  - Improve service metadata endpoints
    - Include available categories

## September 2023 TrustyAI 0.3.0

- *Explainers*
  - Support for external explainability libraries
- *Metrics*
  - Additional metrics
  - Metrics statistical tests

## December 2023 TrustyAI 0.4.0

- *Storage*
  - Wider storage support (database backends)
- *Explainers*
  - NLP explainability support
- *Metrics*
  - Support for user-defined historical windows