*Solution Brief*

# Data Wrangling and Curation for AI

The explosive growth of AI development is facing a hidden roadblock: data preparation. Dirty, inconsistent, and biased data can significantly impact AI model performance and reliability. Furthermore, ensuring data provenance and responsible data handling are crucial for navigating regulatory hurdles and building ethical AI.

**Three key challenges necessitate a robust data wrangling and curation solution for AI development:**

- **The Data Prep Pitfall:** The sheer volume and complexity of data required for AI projects makes data preparation a significant bottleneck. Manually cleaning and organizing data is a time-consuming task that can stall AI development.

- **The Quality Conundrum:** Inconsistent and biased data can lead to unreliable and potentially harmful AI models. Organizations struggle to identify and remove these quality issues and biases within their datasets, compromising the effectiveness and fairness of their AI.

- **The Risk Radar Issue:** Data provenance, or the origin and history of data, becomes a critical concern in the age of AI regulation. Tracing the path of data and ensuring responsible data handling practices are essential for organizations to avoid regulatory issues and ethical concerns.

**40**% of AI adopters struggle with data management aspects like integration, cleaning, and access *(Deloitte)*.

In an age of ever-growing data volumes and complex AI development pipelines, can organizations afford to rely on manual data cleaning processes? Shouldn't AI development be fueled by a reliable and ethical data foundation?

**This is where Zubin steps in.**

Zubin is Data Dynamics' AI-powered self-service data management software, bringing a fresh approach to privacy, security, compliance, governance and optimization in the world of AI-led workloads. It empowers enterprises by enabling users across all levels - from C-suite to data owners - to discover, define, act, transform, and audit data through a user-friendly interface. Zubin brings correlation, consistency and standardization across your organization by delivering granular insights, deriving recommended workflows, and automating actions using personalized policies and RBAC-driven processes. This transformation fosters a culture of data ownership, where everyone becomes a data champion, and the organization fulfills its responsibility as a data custodian.

Zubin addresses these challenges by providing a robust, reliable, and ethical AI data pipeline. Our software tackles data complexity through a suite of features to simplify data preparation tasks, ensuring clean, consistent, and ethically sourced data for AI projects. This leads to faster development, improved model accuracy, and more reliable AI results.

# Integrated Data Acquisition and Consolidation

- **Data Aggregation and Harmonization**
  Aggregates data from diverse sources like file & object storage, and data lakes. Zubin utilizes data mapping and schema management tools to transform data into a unified format, facilitating seamless integration and analysis.

- **Self-Service Data Classifications and Migrations**
  Empowers data owners to classify and migrate their data based on predefined policies. This ensures data is organized and readily available for AI projects.

- **Data Tiering and Placement**
  Accurately identifies, classifies, and maps physical and virtual data locations. Zubin dynamically allocates resources and storage to optimize data accessibility for AI development.

# Automated Data Cleaning and Deduplication

- **Data Classification and Content Analytics**
  Utilizes advanced classification, tagging, and content analysis powered by AI/ML and NLP. This identifies data types, sensitive information (PII/PHI), and potential biases within the data.

- **Data Redundancy Management**
  Detects and eliminates redundant data copies based on metadata discovery, classification, and tagging. This reduces storage requirements and improves data consistency for AI training.

- **Statistical Sampling**
  Selects representative data subsets for efficient metadata and content analysis. This ensures a comprehensive understanding of the data while optimizing processing time.

# Streamlined Data Transformation and Enrichment

- **Data Transformation**
  Facilitates data transformation through various functionalities. Zubin allows data owners and data scientists to clean, transform, and structure data efficiently based on project requirements.

- **Data Archival and Retrieval**
  Enables data archiving in a secure and accessible format (object storage) while facilitating efficient retrieval via software or APIs. This ensures long-term data preservation for potential future AI projects.

- **Data Integration and Harmonization**
  Integrates data from disparate sources and ensures consistency in data format and structure across the entire AI pipeline.

# Enhanced Data Quality and Risk Management

- **Data Observability and Root Cause Analysis**
  Continuously monitors data usage within pipelines using machine learning. This helps identify anomalies, data quality issues, and potential biases, allowing for corrective actions.



**Click here for a demo**

- **Risk Exposure Insights**
  Provides actionable insights using multi-level analytics to identify potential data risks like security vulnerabilities and biases. This enables data scientists to mitigate risks and build trustworthy AI models.

- **Data Security Orchestration and Automation**
  Automates various data security tasks like user activity monitoring and remediation actions based on personalized workflows. This ensures data confidentiality and regulatory compliance throughout the AI lifecycle.

## Fostering Collaboration and Ethical AI Development

- **Data Usage & Traceability**
  Maintains a centralized data index for complete transparency into data usage activities. This facilitates auditability, regulatory compliance, and responsible data handling practices for AI development.

- **Data Ownership Management**
  Ensures clear data ownership throughout the AI pipeline. Zubin empowers data owners to manage access permissions and contribute to data quality checks.

- **Content Analytics for Ethical AI**
  Content analytics powered by AI/ML helps identify potential biases within the data. This empowers data scientists to develop fair and ethical AI models that are free from discrimination.

# What's in it for you?

## Effortlessly Aggregate and Harmonize Data

Studies show 76% of individuals want more control over their data. Limited data visibility creates blind spots for privacy risks. Zubin offers enterprise-wide risk analysis, providing a consolidated view for proactive mitigation and fostering a culture of privacy by design.

## Automate Data Prep with AI-Powered Cleaning

Research indicates organizations using data quality automation see a 25% reduction in data preparation time for AI. Zubin leverages AI to identify data types, biases, and redundancies. This automates data cleaning and deduplication, saving valuable time for data scientists.

## Empower Data Owners for Self-Service Management

Studies reveal empowered data owners drive a 15% increase in AI project completion rates. Zubin empowers data owners to classify and migrate data based on policies, fostering data ownership, improving organization, and ensuring readily available data for AI development.

## Build Reliable AI Models with Enhanced Data Quality

Research suggests robust data quality practices in AI development lead to 30% improvement in model accuracy. Zubin offers data observability, risk exposure insights, and data security automation. Enterprises can easily identify data issues, biases, and security risks, and thereby empower data scientists to build reliable and trustworthy AI models.

## Promote Collaboration and Ethical AI Development

Studies show prioritizing responsible AI practices increases public trust by 20%. Zubin provides data usage & traceability, data ownership management, and content analytics for ethical AI. This fosters clear data ownership, transparent data usage, and empowers data scientists to develop fair and unbiased AI models.

Your next chapter of success awaits; let's write it together with Zubin.

**Click here for a demo**