# End-to-end Integrated Platform for Semiconductor Analytics



**Information Security and IP Controls**

**AI/ML**
- Guided Analytics
- Advanced Insights for Manufacturing
- Edge Analytics
- ModelOps

**Data Intelligence**
- End-to-end traceability
- Feature extraction
- Semiconductor specific data models

**Control**
- Process Control
- Test Control
- Assembly Control

**Integration Services**
- Sapience Manufacturing Hub
- Data Exchange Network (DEX)

**IoT Software**
- PDF Created Data
  - Design Characteristics
  - In-Field Data
  - pdFasTest
  - eProbe
- Customer Data
  - Process & Eq. Data
  - Fab Data
  - Assembly Data
  - Test Data
- Fab Equipment and Fabless Integrations based on Standards

**APIs**

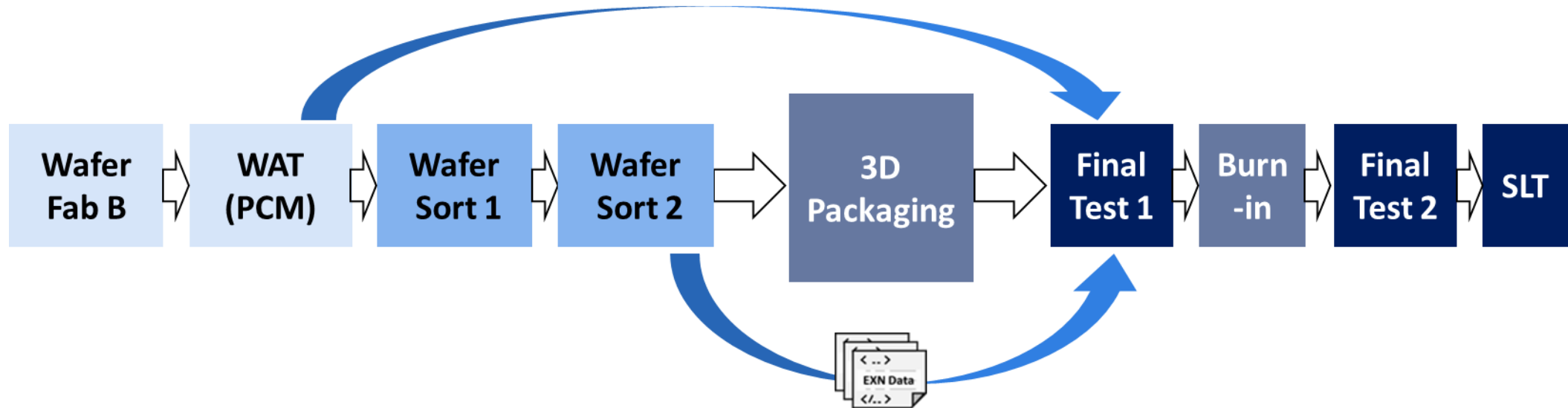**Enterprise Integrations (Data Lakes, ERP, PLM, CRM...)**

**Fully integrated solution to accelerate production ramp, improve overall yield and quality for Semiconductors**

3

# Overview

- **Application Scenario**

- **Edge Deployment Architecture**

- **Model Training and Distribution**

- **Exensio Data Feed Forward Integration**

- **Demo: Data Feed Forward Orchestration**

- **Demo: Edge Model Execution**
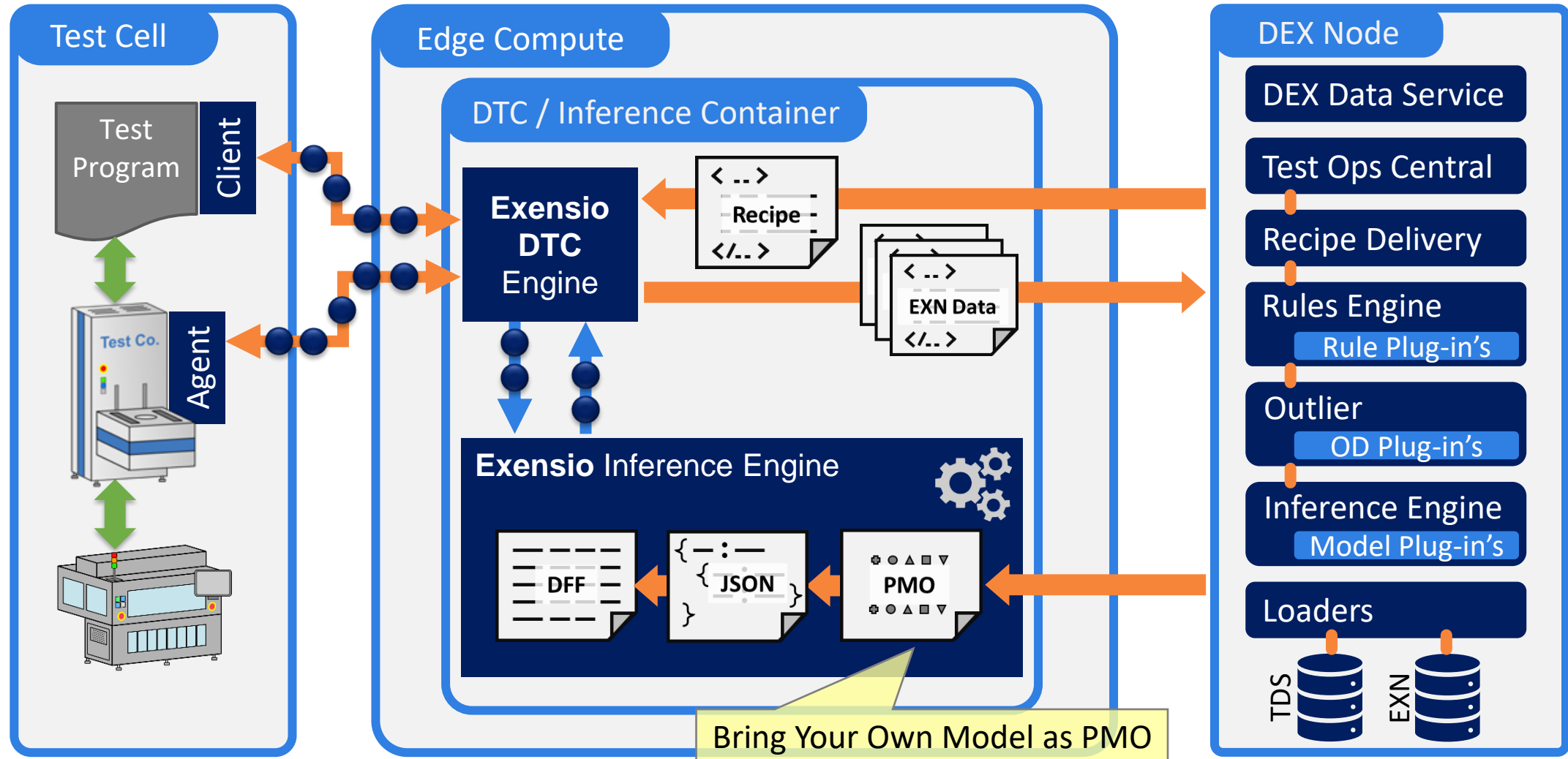
# Backdrop Scenario



- **Deploy a feed-forward-loop test operation for a multi-chip device**
  - Reliable and timely capture of wafer sort test operation data
  - Automate feature extraction from PCM/WAT & WS to produce feed-forward data
  - Mechanism to train & periodically retrain models on current data
  - Deploy models across supply chain to the package test operation
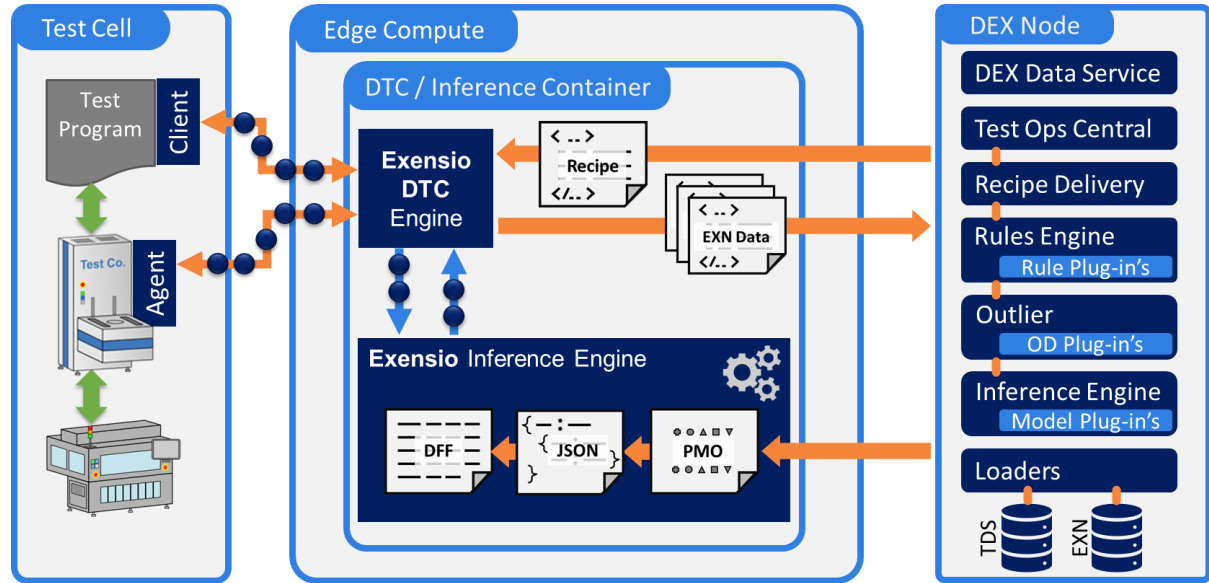  - Monitor test data and model performance

# Edge Integrated Inference Container

*Your model with Exensio rules & model management deployed to the edge*



Bring Your Own Model as PMO

PDF/SOLUTIONS™

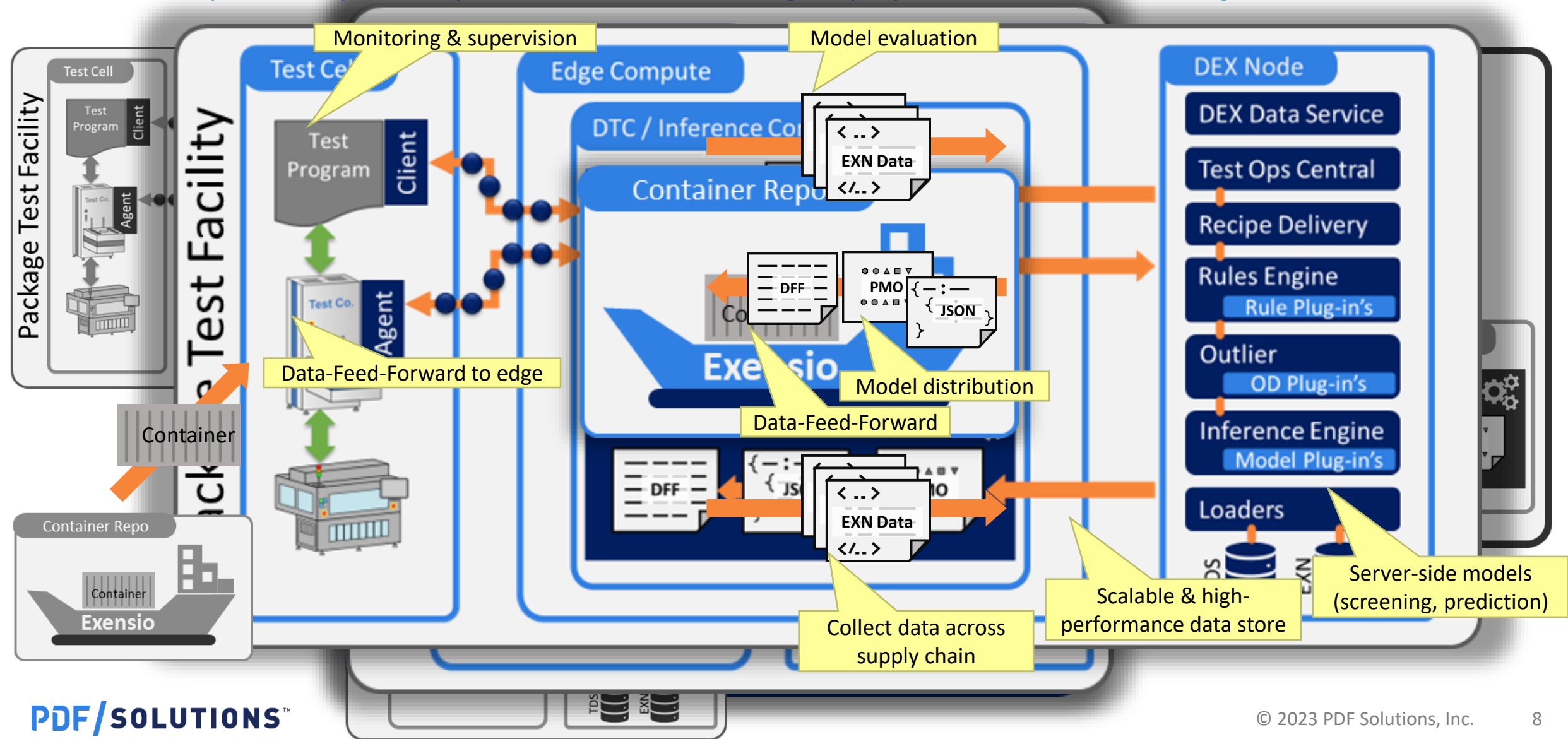# Edge Integrated Inference Container

*Your model with Exensio rules & model management deployed to the edge*



- Synchronous inline inference per test flow enables adaptive test

- High-speed prediction and bin override (<200ms roundtrip)

- At-scale deployment architecture

- Bring Your Own (BYO) model

- Full spectrum data feed

- Compliment model with rules
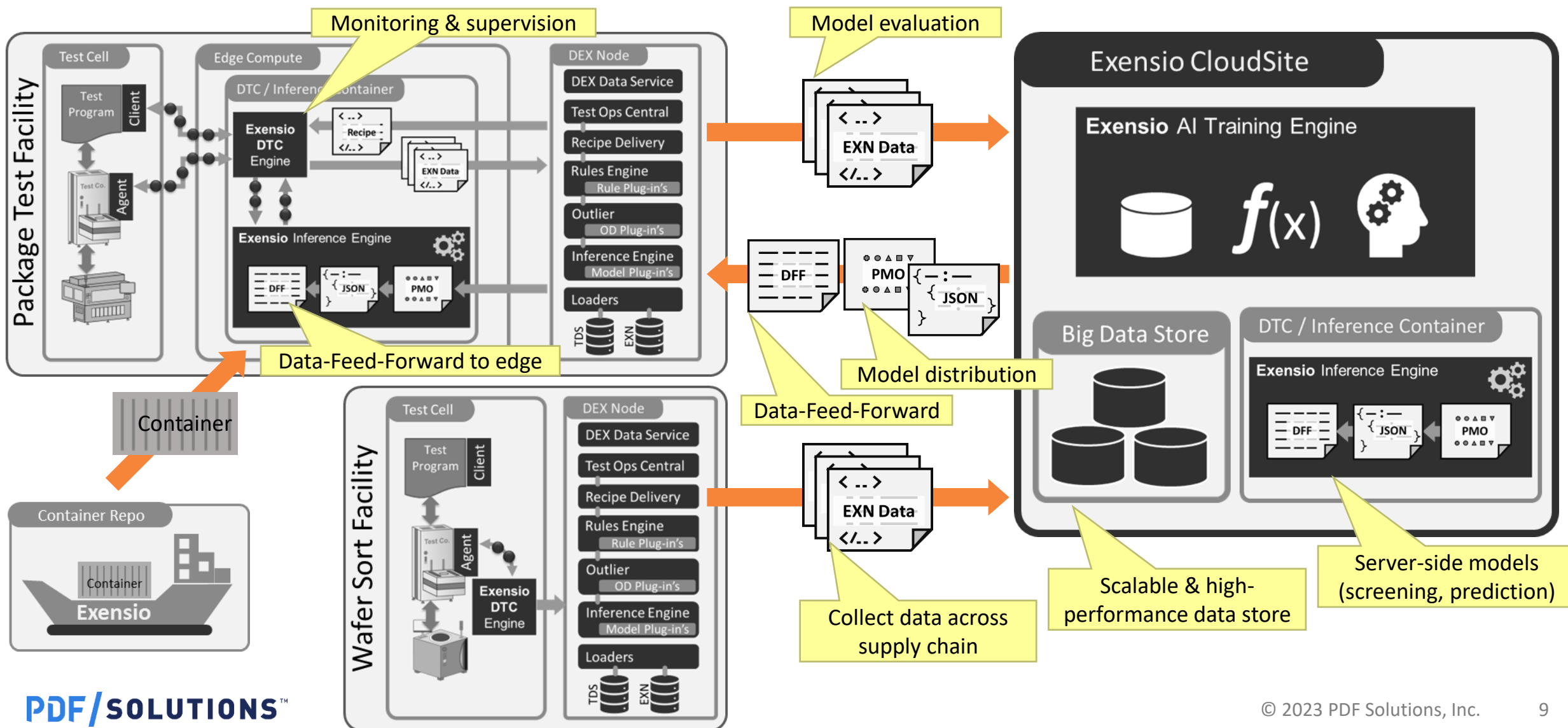
- Secure execution environment

# Exensio ML Model Deployment

*Automatic dynamically scaled prediction model training, deployment and model management*
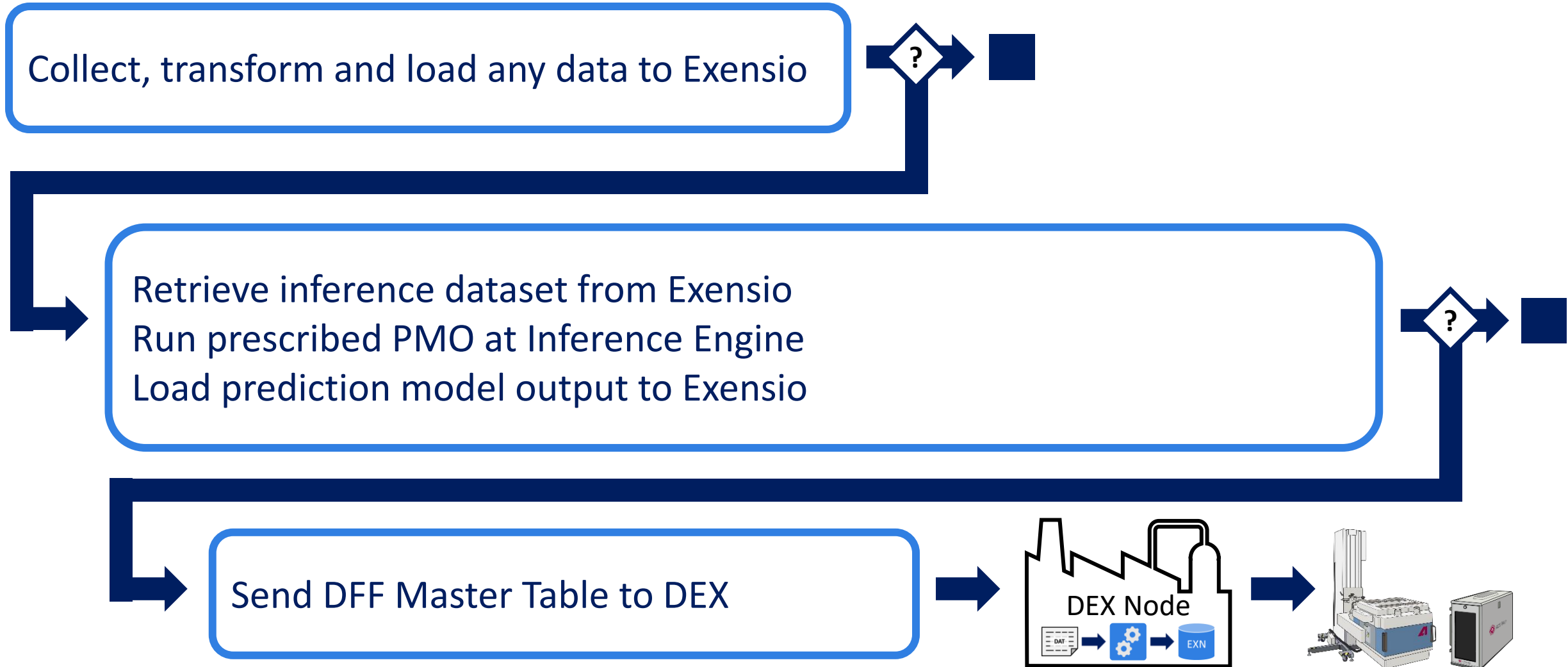


Monitoring & supervision

Model evaluation

Data-Feed-Forward to edge

Model distribution

Data-Feed-Forward

Collect data across supply chain

Scalable & high-performance data store

Server-side models (screening, prediction)

Test Cell

Edge Compute

DTC / Inference Control

Container Repo

EXN Data

DFF

PMO

JSON

Exensio

DEX Node

DEX Data Service

Test Ops Central

Recipe Delivery

Rules Engine
Rule Plug-in's

Outlier
OD Plug-in's

Inference Engine
Model Plug-in's

Loaders

Package Test Facility

Test Cell

Test Program

Client

Test Co.

Agent

Container

Container Repo

Container
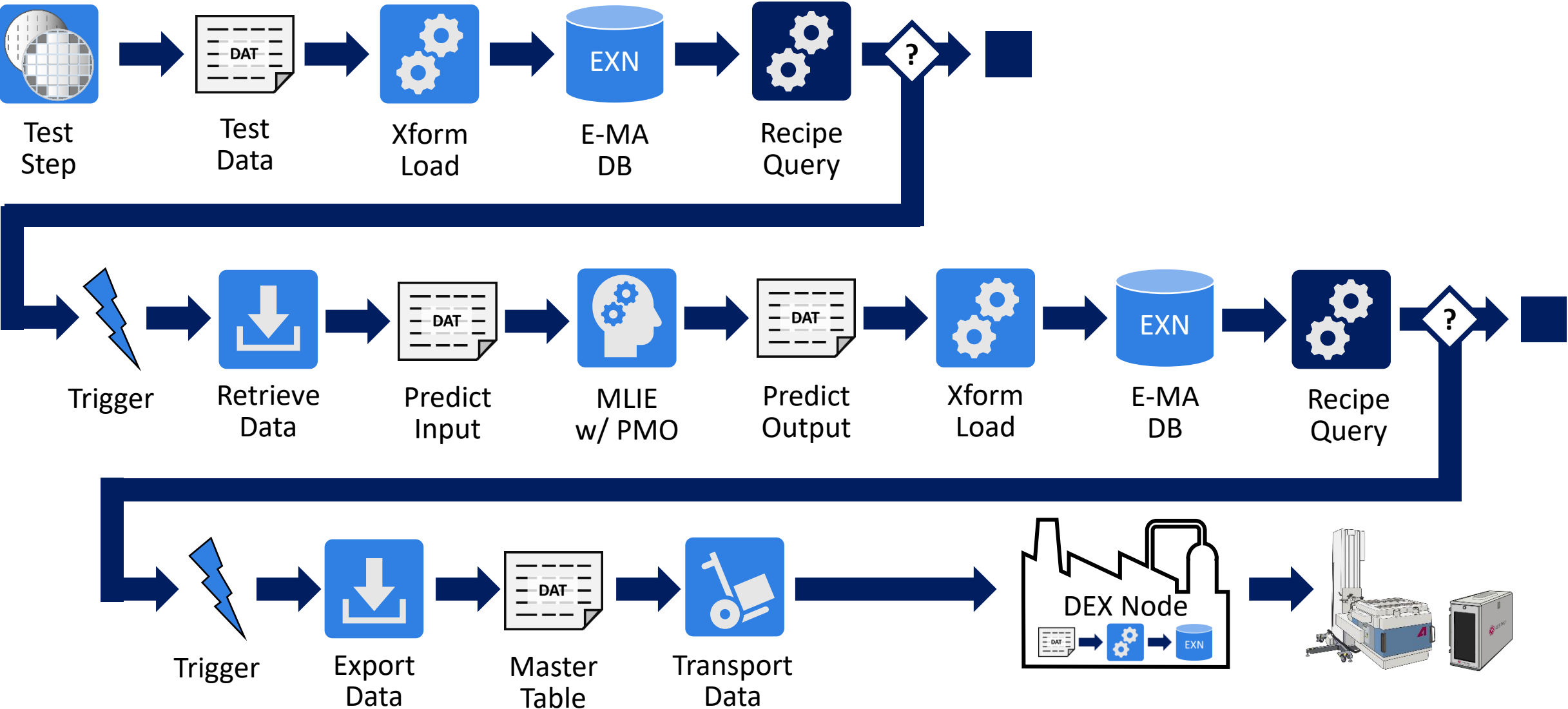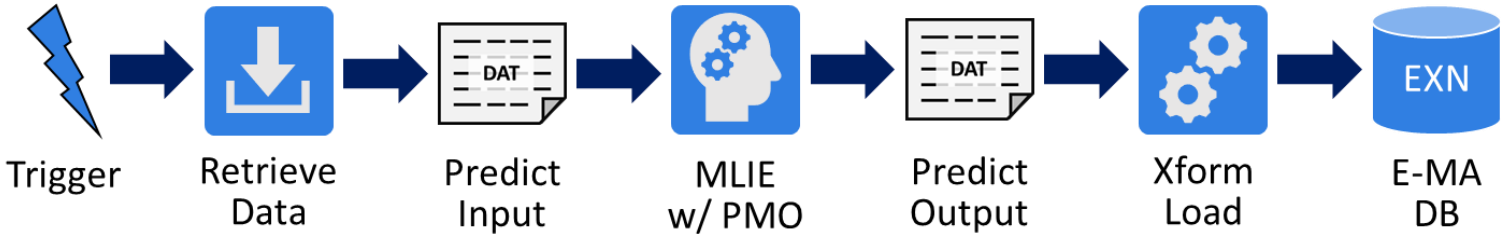
Exensio

PDF/SOLUTIONS™

# Exensio ML Model Deployment

*Automatic dynamically scaled prediction model training, deployment and model management*

# ML Inference Engine Prediction and Data-Feed-Forward

Collect, transform and load any data to Exensio

Retrieve inference dataset from Exensio
Run prescribed PMO at Inference Engine
Load prediction model output to Exensio

Send DFF Master Table to DEX

DEX Node

# ML Inference Engine Prediction and Data-Feed-Forward



Test Step → Test Data → Xform Load → E-MA DB (EXN) → Recipe Query → ? →

Trigger → Retrieve Data → Predict Input → MLIE w/ PMO → Predict Output → Xform Load → E-MA DB (EXN) → Recipe Query → ? →

Trigger → Export Data → Master Table → Transport Data → DEX Node →

# Per-Wafer DFF Model Outputs

Trigger → Retrieve Data → Predict Input → MLIE w/ PMO → Predict Output → Xform Load → E-MA DB

**Wafer-level statistics for parameters used in die-level prediction model**

| MLIE_Wafer_Master_Table | | | | | |
|---|---|---|---|---|---|
| Lot | Wafer | start_time | PCM_PROG_002.Ib_p | PCM_PROG_002.Vs_p | PCM_PROG_002.Is_n |
| LOTID1 | 1 | 5/4/2023 2:14:18 PM | 0.16 | 0.12 | 0.19 |
| LOTID1 | 2 | 5/4/2023 2:14:59 PM | 0.08 | 0.64 | 0.11 |
| LOTID1 | 4 | 5/4/2023 2:15:35 PM | 0.81 | -0.21 | 0.00 |
| LOTID1 | 5 | 5/4/2023 2:15:50 PM | 0.65 | -0.34 | 0.01 |
| LOTID2 | 1 | 5/4/2023 2:14:40 PM | -0.27 | -0.41 | 0.38 |
| LOTID2 | 2 | 5/4/2023 2:14:43 PM | 0.19 | -0.30 | -0.34 |
| LOTID2 | 3 | 5/4/2023 2:14:43 PM | -0.10 | -0.13 | -0.36 |
| LOTID2 | 4 | 5/4/2023 2:15:50 PM | 0.06 | 0.40 | 0.74 |
| LOTID2 | 5 | 5/4/2023 2:15:50 PM | 0.10 | 0.08 | -0.73 |
| LOTID3 | 1 | 5/4/2023 2:16:06 PM | -0.23 | -0.47 | -0.30 |
| LOTID3 | 2 | 5/4/2023 2:16:09 PM | 0.20 | 1.66 | 0.41 |

# Per-Die DFF Model Outputs



Trigger → Retrieve Data → Predict Input → MLIE w/ PMO → Predict Output → Xform Load → E-MA DB

**Die-level engineered feature from raw DFF wafer sort parametric inputs**

**MLIE_Wafer_Master_Table**

| Lot | Wafer | ecid | WS_PROG_002.simple_Score | WS_PROG_002.BVces_1 | WS_PROG_002.IDDQ_2 | WS_PROG_002.VDDmin_2 | WS_PROG_002.LKG_2 | WS_PROG_002.BVces_2 |
|---|---|---|---|---|---|---|---|---|
| LOTID1 | 1 | Lotid1_1_-14_-2 | 0.22 | 0.21 | -0.16 | -0.07 | 0.47 | 0.23 |
| LOTID1 | 1 | Lotid1_1_-14_0 | 0.22 | 0.38 | 0.14 | 0.34 | 0.45 | 0.24 |
| LOTID1 | 1 | Lotid1_1_-13_1 | 0.23 | 0.30 | 0.31 | -0.21 | 0.21 | 0.31 |
| LOTID1 | 1 | Lotid1_1_-13_3 | 0.22 | 0.20 | 0.23 | 0.21 | 0.18 | 0.24 |
| LOTID1 | 1 | Lotid1_1_-12_-4 | 0.22 | 0.28 | 0.38 | 0.60 | 0.29 | 0.16 |
| LOTID1 | 1 | Lotid1_1_-12_-2 | 0.22 | 0.28 | 0.16 | 0.11 | 0.25 | 0.22 |
| LOTID1 | 1 | Lotid1_1_-12_0 | 0.22 | 0.27 | 0.31 | 0.32 | 0.23 | 0.23 |
| LOTID1 | 1 | Lotid1_1_-12_1 | 0.23 | 0.28 | 0.35 | 0.00 | 0.26 | 0.30 |
| LOTID1 | 1 | Lotid1_1_-12_2 | 0.22 | 0.33 | 0.11 | 0.54 | 0.31 | 0.32 |
| LOTID1 | 1 | Lotid1_1_-12_4 | 0.22 | 0.22 | 0.02 | -0.31 | 0.25 | 0.28 |
| LOTID1 | 1 | Lotid1_1_-11_-5 | 0.22 | 0.29 | 0.17 | 0.49 | 0.25 | 0.24 |
| LOTID1 | 1 | Lotid1_1_-11_-3 | 0.22 | 0.27 | 0.22 | 0.05 | 0.30 | 0.26 |

**Filtered subset of raw wafer sort parameters used by prediction model**
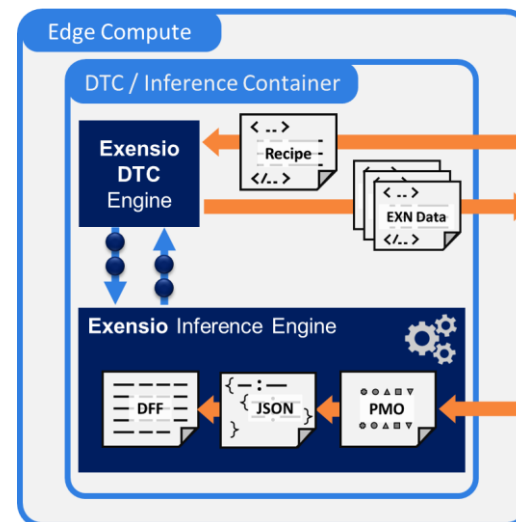
# Recipe Delivery Web Service Model Integration



realtime_1.pmo
and
realtime_1.json
included and delivered within the test session recipe and sourced to Inference Engine container
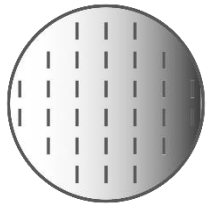
# Edge Data Feed Forward Query Web API

localhost/DataRetrieval/v1/data_feed_forward/query?LotId=XGB017&WaferId=17

{"data":[{"lotId":"XGB017","waferId":"17","program":"F891DA0978175667E0536E30200A9DF8_xgboost","lgKey":21,"param
["wf_key","ecid","WS_PROG_001.Main.analogPwrUp2_A.getVmonPllDistFmt.Main.analogPwrUp2_A.getVmonPllDistFmt_t0_Pad
"rows":[
{"values":[2,"XGB017_17_-7_-10",0.51,0.5,0.5,0,0.2989977]},
{"values":[2,"XGB017_17_-7_-8",0.48,0.48,0.48,0,0.040815283]},
{"values":[2,"XGB017_17_-7_-6",0.47,0.48,0.48,0,0.040815283]},
{"values":[2,"XGB017_17_-6_-12",0.5,0.5,0.51,0,0.2989977]},
{"values":[2,"XGB017_17_-6_-7",0.48,0.48,0.48,0,0.040815283]},
{"values":[2,"XGB017_17_-6_-5",0.47,0.47,0.47,0,0.040815283]},
{"values":[2,"XGB017_17_-6_-3",0.47,0.48,0.47,0,0.040815283]},
{"values":[2,"XGB017_17_-6_10",0.49,0.49,0.5,0,0.2989977]},
{"values":[2,"XGB017_17_-6_12",0.51,0.52,0.5,0,0.4607918]},
{"values":[2,"XGB017_17_-5_-11",0.5,0.5,0.5,0,0.2989977]},
{"values":[2,"XGB017_17_-5_-9",0.48,0.49,0.48,0,0.07394859]},
{"values":[2,"XGB017_17_-5_-6",0.48,0.47,0.47,0,0.040815283]},
{"values":[2,"XGB017_17_-5_0",0.47,0.47,0.47,0,0.040815283]},
{"values":[2,"XGB017_17_-5_2",0.47,0.47,0.47,0,0.040815283]},
{"values":[2,"XGB017_17_-5_8",0.48,0.48,0.48,0,0.040815283]},
{"values":[2,"XGB017_17_-5_11",0.51,0.5,0.5,0,0.2989977]},
{"values":[2,"XGB017_17_-4_-12",0.5,0.5,0.5,0,0.2989977]},
{"values":[2,"XGB017_17_-4_-7",0.47,...,0.4...,0,0.040815283]},

**Data Feed Forward** dataset retrieved from web service API with array of multiple parameters **per-device** by ECID and sourced to Inference Engine container
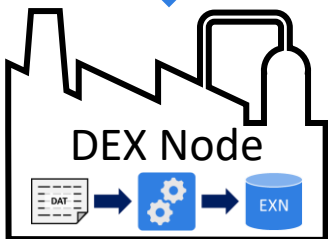
# Demo Scenario

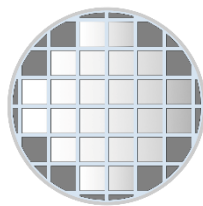| Define PCM to Wafer-Level DFF Transform Rule | Define Wafer Sort to Die-Level DFF Transform Rule | Define ML Edge Prediction Rule | Execute ML Model Online / Inline at Final Test |
|---|---|---|---|
|  |  |  | DEX Node |
| Define Rule Scope | Define Rule Scope | Define Rule Scope | Test Program Init |
| Upload PMO / JSON to Transform Data & Engineer Features | Upload PMO / JSON to Transform Data & Engineer Features | Upload Prediction Model as PMO / JSON | Container Launch |
| | | Define Real Time Action | Retrieve PMO / JSON |
| | | | Real Time Predictions |
| DEX Node | DEX Node | DEX Node | |

PDF/SOLUTIONS™

# Demo Videos



MLIE Data-Feed-Forward Workflow — Recorded Demonstration, May 2023

03:51



Exensio Inference Container for ACS Edge — Recorded Demonstration, August 23, 2022

04:50 – 07:54

# Statistical and ML Apps for ACS Edge™

## Statistical Applications

- **ACS Outlier Screening:** Real-time outlier screening using a variety of algorithms

- **ACS Adaptive Test:** Real-time adaptive test; Test More or Test Less

- **ACS Tester Control:** Real-time statistical process control (SPC) to avoid quality excursions and escapes

- **ACS Statistical Binning:** Dynamically bin devices via statistical rules and without test program modification

## ML Based Applications

- **ACS Predictive Binning:** Dynamically bin devices with upstream data (Data Feed Forward) and via an ML model decision

- **ACS Outlier Screening ML:** Real-time outlier screening using ML

- **ACS Adaptive Test ML:** Real-time adaptive test using ML

- **ACS Custom ML:** Bring-Your-Own ML model and leverage the Exensio Data Feed Forward/Backward infrastructure