# Galaxy Data Quality Program
# MIT IQ Industry Symposium
## July 18-19, 2007

Ingenix
United Health Analytics
Galaxy – Shared Data Warehouse
Laura Sebastian-Coleman
IS Manager – Data Quality & End User Support

**INGENIX** ®

# Overview

- Ingenix and Galaxy
- Galaxy's DQ program
- Evolving business needs and the pace of change
- Data quality in relation to evolving business needs
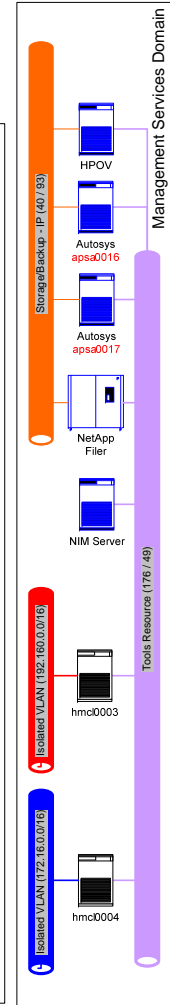
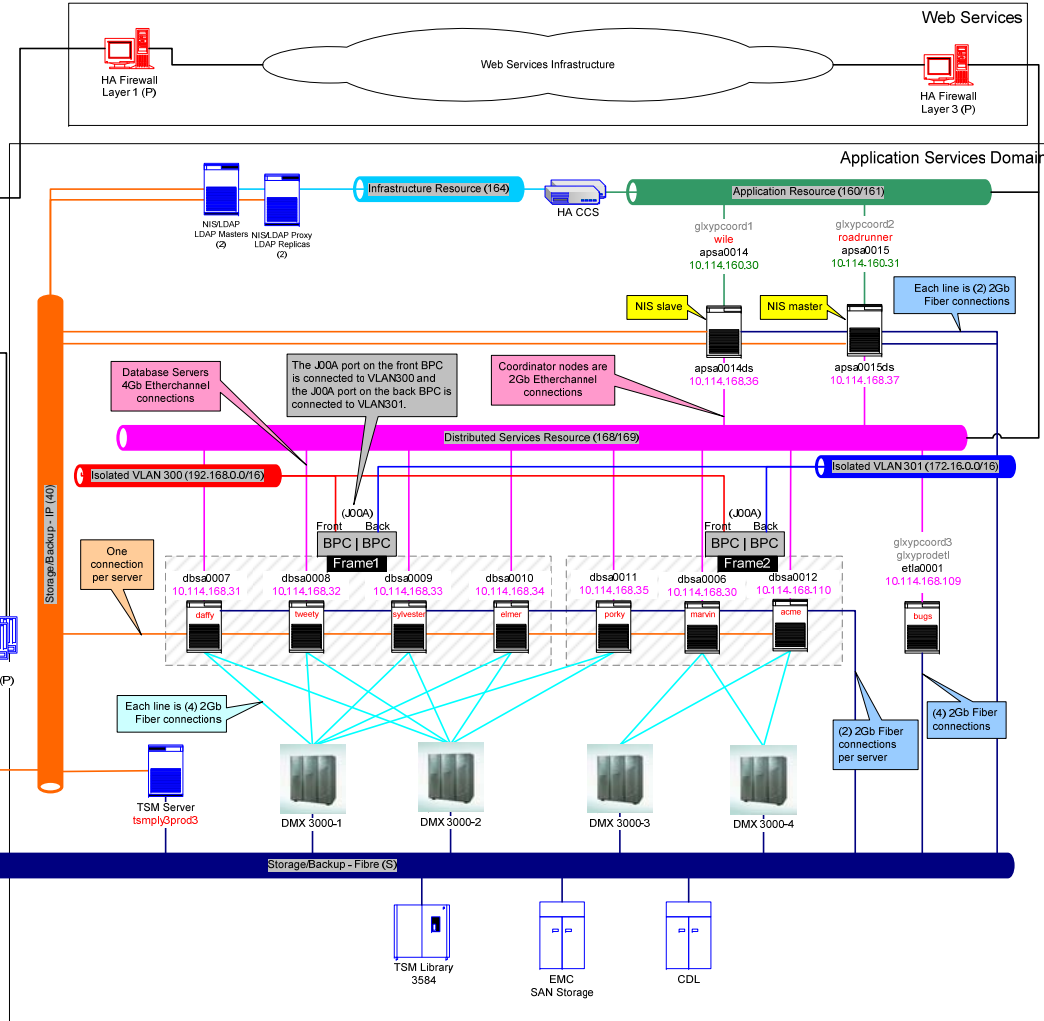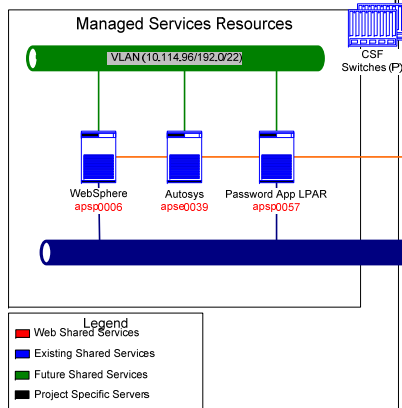**INGENIX.**

# Ingenix Background

- A global healthcare information company
- Founded in 1996 to develop, acquire, and integrate some of the nation's best-in-class healthcare information capabilities
- Significant and rapidly evolving portfolio of tools and services now transform data into actionable, fact-based, technology-enabled decision support
- Ranked among the top 10 providers of informatics by *Healthcare Informatics* magazine in June 2006
- Today there is an Ingenix product at work in nearly every U.S. healthcare organization.
- Ingenix is a wholly owned subsidiary of UnitedHealth Group (UHG).

**INGENIX**®

# Galaxy Overview

- Atomic Data Warehouse with transformations
- Integrates data from more than a dozen subject areas (claim, membership, customer, provider, etc.) across multiple sources
- Size
  - 350 source input files from more than 25 distinct internal and external sources (and counting)
  - 18 TB of data; 62 TB footprint
  - 3,159 attributes across 12,632 columns in 600 tables (and counting)
  - Largest table: more than 1.5 billion rows
    - 1,704,717,031 on Claim Statistical Service as of 5/3/07
- Usage
  - Over 1,000 registered users
  - 7,888 queries per day / 256,656 per month, on average
  - Ad hoc, scheduled queries, production extracts to applications and marts
  - Direct access to Galaxy via user-selected tools – Sagent is administratively supported

**INGENIX.**

# Galaxy Physical Architecture

# System Components

- Hardware
  - 7 IBM P-series Servers P575
  - 2 IBM P-series Servers P510
  - 1 IBM P-series Server P570
  - 4 EMC DMX 3000 Storage Cabinets
  - Additional supporting servers for Sagent, Autosys, etc.
- Software
  - UDB with DPF v8.2
  - AIX 5.3.0
  - DataStage/PX 7.0.1
  - Optiload 3.1
  - CoSort 7.5.3
  - Autosys 4.5
  - Sagent 4.5i

**INGENIX.**

# Galaxy Source Systems & Subject Areas

| Source System | Description | Subject Area | Provides |
|---|---|---|---|
| **ACIS** | Feeds Customer data | | |
| **PRIME** | Feeds Customer Small Group data | | |
| **COSMOS** | Feeds Customer data | **Customer** | Provides Customer demographic, policy, and coverage information |
| **MAMSI LIVE!** | Feeds Customer data | | |
| **Sales Area** | Feeds Customer Sales data | | |
| **Entity Sales** | Feeds Customer data | | |
| **CES** | Feeds Member, Customer coverage data | | |
| **COSMOS** | Feeds Member data | | |
| **MAMSI LIVE!** | Feeds Member data | **Member** | Provides Member demographics and coverage information. |
| **Conversion Notification File** | Feeds UNet Member& Customer data | | |
| **URSULA** | Feeds Individual ID and other identifiers | | |
| **NDB** | Feeds Provider Contract & specialty code data | **Provider** | Provides Provider demographic and contract information |
| **Medco** | Feeds Pharmacy Claim data | | |
| **NCPDP** | Feeds Pharmacy Provider data | **Pharmacy** | Provides detail information from Managed Pharmacy Vendors |
| **First Data Bank** | Feeds NDC Drug and Drug Pricing data | | |
| **TOPS** | Feeds Statistical Claim data | **Statistical Claim** | Provides detail medical claim and service information |
| **COSMOS** | Feeds Statistical Claim data | **Statistical Claim Aggregations** | Provides Inpatient Confinemen and Outpatient Event statistical information |
| **MAMSI LIVE!** | Feeds Statistical Claim data | | |
| **HSS** | Feeds DRG Grouper data | | |
| **UCAS** | Feeds Financial Claim data | **Financial Claim** | Provides detail financial transaction information. |
| **Lab Corp** | Supplies universal names & codes for identifying lab & clinical lab test results. | **Lab Results** | Provides data about performed lab tests & associated results |
| **QUEST** | Feeds Lab Results data | | |
| **USPS File** | Feeds Zip Code data | **Geography** | Provides ZIP and HCFA County Code information |
| **NSAA Market File** | Feeds Product Service Area and Organization data | **Product** | Provides hierarchical product grouping/service area information |
| **Financial Systems** | Feeds Financial market data | **Organization** | Provides Site, Legal Entity, and Market information. |

©2006 Ingenix, Inc

**INGENIX.**

# Functions of Galaxy Data

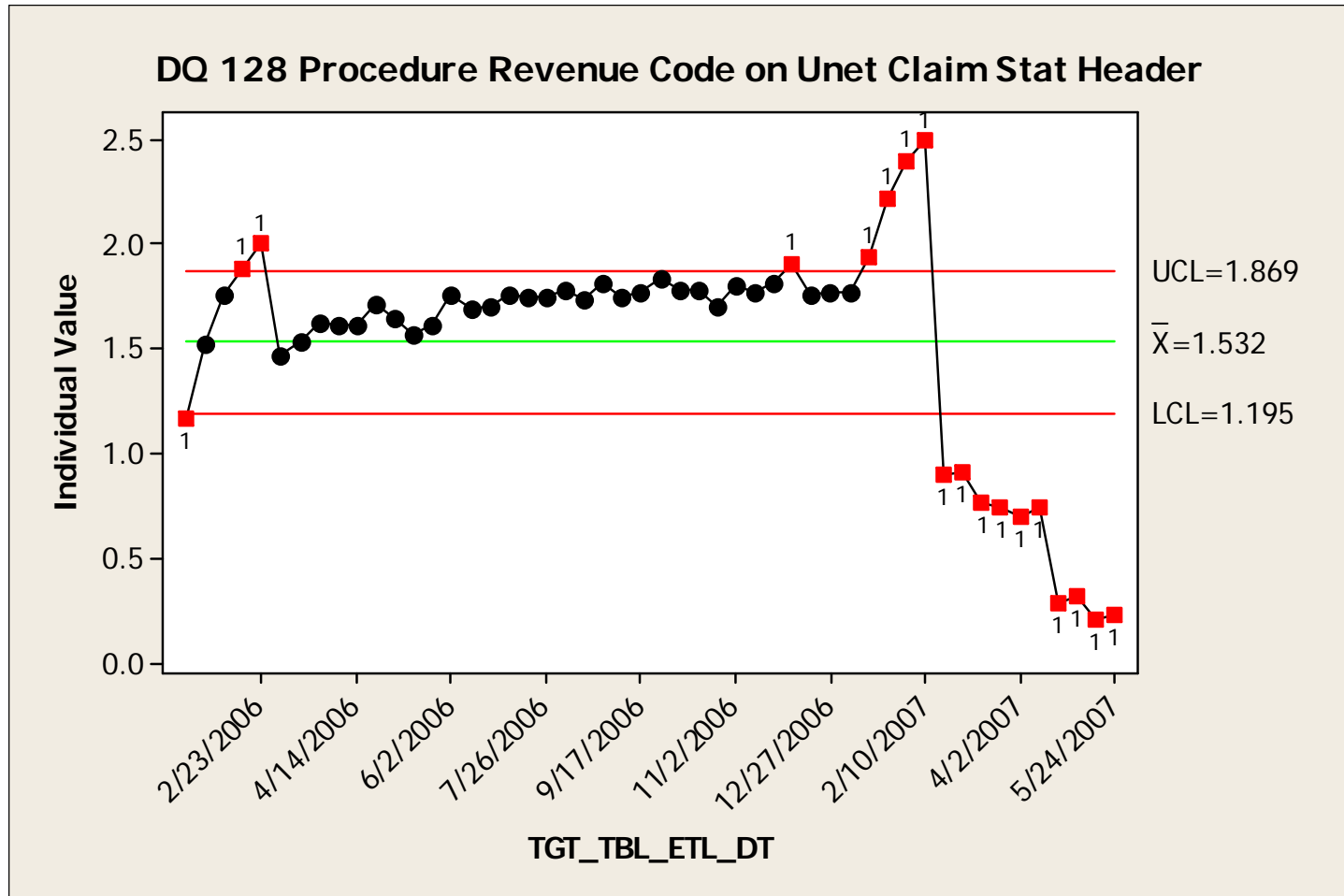Galaxy is the single source of truth for key business functions

- Medical Trend Analytics
- Pricing
- Provider Utilization & Profiling
- Appropriateness of Care
- Network Adequacy
- Care Management / Pattern of Care / Preventive Care
- Fraud & Abuse
- Customer Reporting
- HEDIS Reporting
- Member Demographics
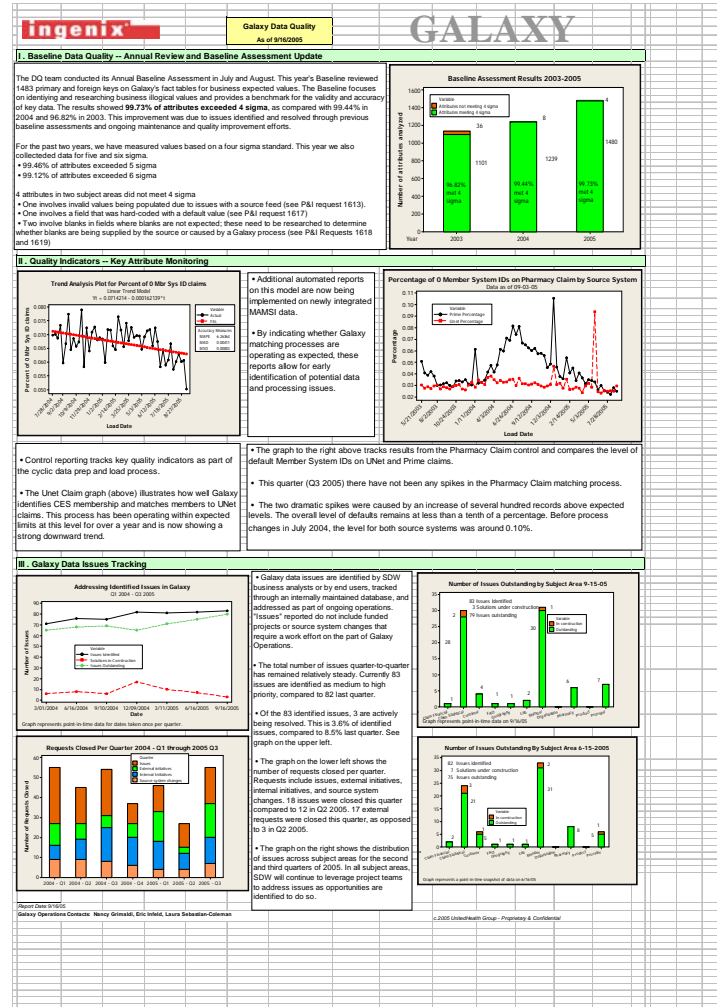- Product Penetration

**INGENIX**®

# Galaxy's Data Quality Program

- Management recognized need for DQ when Galaxy was launched
- Theoretical / methodological foundations
    - Correct data problems at the source
    - Data as a product
    - Statistical process control
- Primary functions of DQ program
    - Monitor, measure, and report on Galaxy's Data Quality
    - Recommend and implement actions based on findings
- Biggest initial challenge = establishing useful metrics
    - What to measure / how to measure
    - How to respond to the results of measurements
- 2003 Initiated metrics & reporting program
- 2004 Implemented first automated measures
- 2004-2007: Deliver weekly/cyclic, monthly, quarterly, semi- annual reporting through largely automated processes

**INGENIX.**

# Example of Weekly Measure



DQ 128 Procedure Revenue Code on Unet Claim Stat Header

INGENIX.

# Quarterly Management Report

- Baseline Data Quality – annual review and baseline assessment

- Quality Indicators – key attribute monitoring

- Galaxy Data Issues Tracking

**INGENIX**®

# Current Situation

- Galaxy = a mature, enterprise data warehouse
- High demand for data and for organizational services
- Galaxy's DQ program also relatively mature
    - Defined metrics
    - Automated data collection
    - Regular reporting
    - DQ Community
- UHG growing, largely through acquisitions and partnerships
- Healthcare industry changing – relation of government to health care, new products, esp. consumer driven

**INGENIX.**

# Pace of Change for Galaxy

- 2004
  - Galaxy integrated data from MAMSI, a United Health Group acquisition
    - Used the existing structure
    - 1+ year to integrate
- 2006
  - Integrated data from three new source systems
  - Developed a new subject area, Revenue
  - Significantly expanded Customer subject area
  - Responded to healthcare industry changes
    - Part D data
    - HRA (Health Reimbursement Account) data
- 2007
  - Integrate data from additional acquisitions
  - Expand the Revenue subject area
  - Continue to support the use and enhancement of existing data.
- 2008
  - Two major integrations already scheduled
  - Potential for several others

**INGENIX®**

# Pace of Change for Galaxy DQ

- Biggest challenge
    - 2003 what to measure and how to measure
    - 2007 how to rapidly analyze and act on DQ data
- Baseline Assessment of Galaxy Data Quality
    - 2003
        - 800 person hours to pull and analyze data for first Baseline Assessment
        - Duration = more than 3 months
        - Measured 1137 attributes
    - 2006
        - Pulled 75% of data in less than 10 hours through an automated process
        - Measured 1506 attributes
        - Pull data quarterly
- Automated reports
    - 2004:  4 reports
    - 2007: 80 reports
    - Reports now implemented as part of standard development process.

**INGENIX.**

# 2007 – 2008 Key UHG Business Needs

- UHG acquisitions and partnerships –
    - More data for Galaxy
    - More users need access
- Users need data sooner –
    - Time to integrate data into Galaxy must be shortened
- Legacy data critical for ensuring reporting continuity and analytics –
    - Continued support is necessary
- Data consistency across sources critical for reporting continuity and analytics –
    - Integration methodologies need to promote and enforce consistency

**INGENIX.**

# How to Respond?

- Data Quality included in set of changes to improve efficiency and agility
  - Common Interface – puts more responsibility on source systems for data quality
  - Gateway – changes how Galaxy prepares data.
- DQ measures
  - More comprehensive
  - Taken earlier in the process
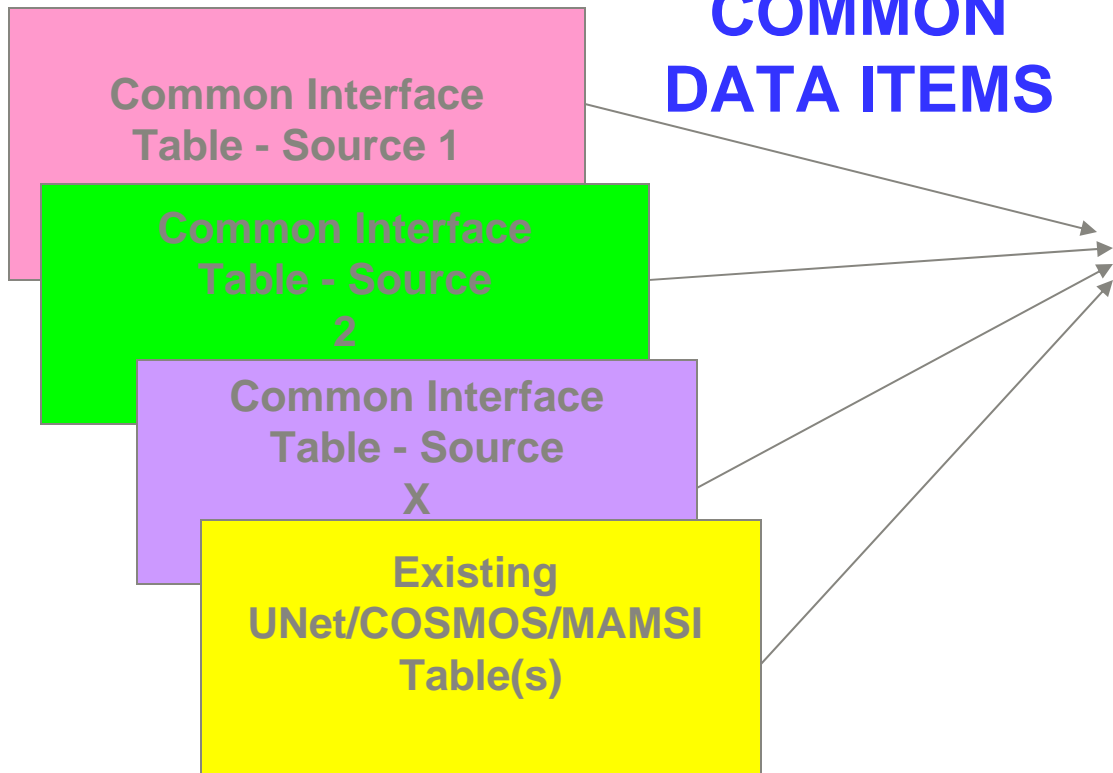  - More fully automated
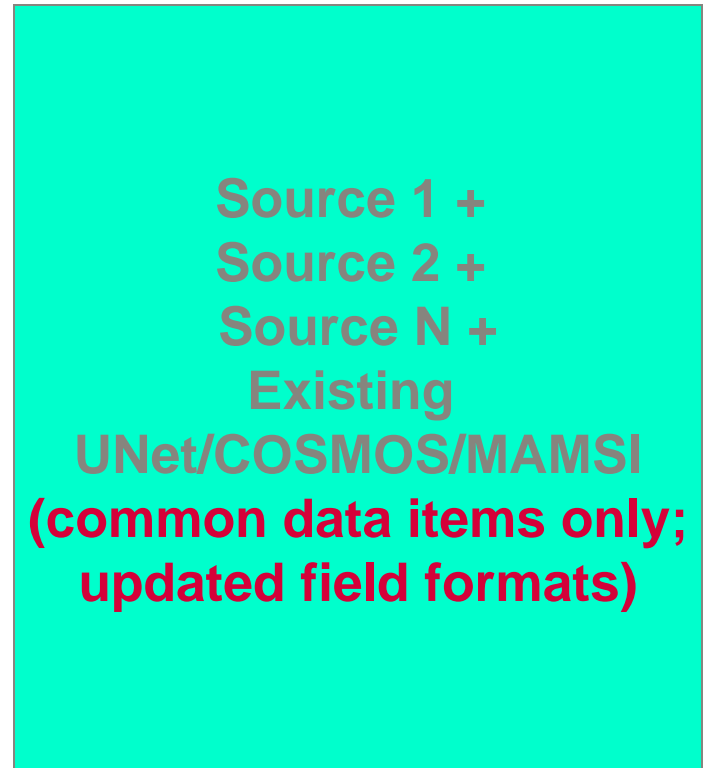
**INGENIX.**

# Common Interface Approach

- Galaxy defines standard requirements and layouts for data
- Sources map to these requirements and feed to Galaxy
- Streamlined transformation/load into Galaxy
- Common model across the enterprise

**INGENIX.**

# Common Interface Architecture – Views

**Physical Tables (Objects)**

**Enterprise View**

**COMMON DATA ITEMS**

Common Interface Table - Source 1

Common Interface Table - Source 2

Common Interface Table - Source X

Existing UNet/COSMOS/MAMSI Table(s)

Source 1 +
Source 2 +
Source N +
Existing
UNet/COSMOS/MAMSI
**(common data items only; updated field formats)**
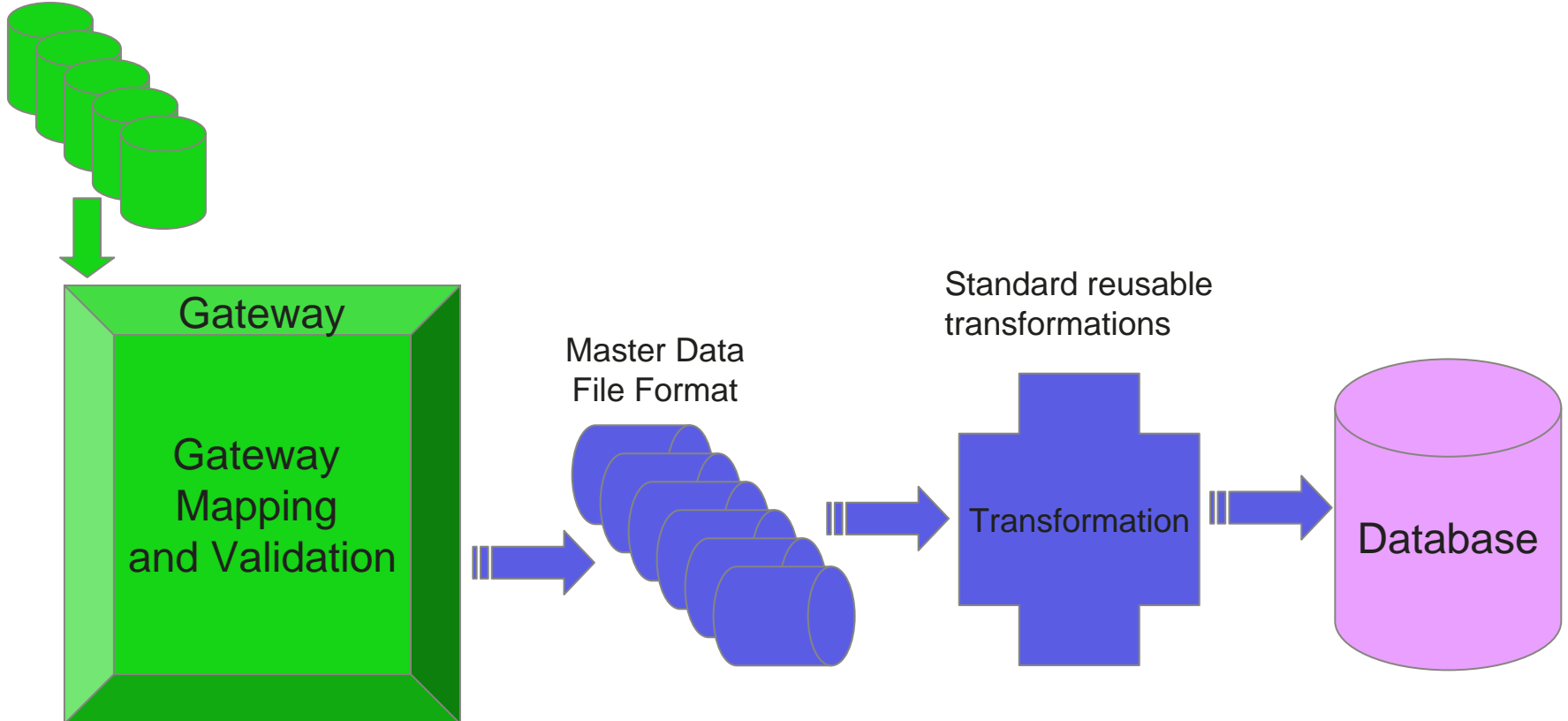
**INGENIX.**

# Gateway Integration Tool

- Facilitates mapping disparate data sources into a Master Data Definition
- Applies generic transformation logic to the output
- Utilizes reusable transforms
- Performs automatic code generation
- Ensures consistency across source-to-target mappings
- Provides true-to-code documentation
- Incorporates data quality modules
- Increases speed and reduces complexity of data integrations

**INGENIX.**

# Gateway Integration Tool

New Data
Sources

Gateway

Gateway
Mapping
and Validation

Master Data
File Format

Standard reusable
transformations

Transformation

Database

**INGENIX.**

# Gateway – Data Quality Features

- DQ functions
    - Monitor and react to events in processing
    - Collect trend data
- Field validation
    - Data type checking
    - Value range checking
    - Valid value list checking
    - Assignment of default values
    - Informational, error and warning messages
- File validation
    - Format checking
    - Field counts / record length validation
    - Summary of field error and warning messages
    - Thresholds of summary counts of errors and warnings that allow job to be aborted if counts or percentages exceeded – generate alerts

**INGENIX.**

**INGENIX**®

# Back to Basic DQ

- Data in the warehouse is only as good as data in the source
    - Ensuring sources to supply better data through the Common Interface
- Manufacturing model: Data as a product produced through a process
    - Executing processes more consistently across the database through the Gateway
- Measure to improve
    - Gateway integrates and executes DQ measures consistently across the database.
    - Both tools measure ETL processes (timing of jobs, etc.) that affect other aspects of data quality from end-to-end

**INGENIX**®

# DQ: Chicken or Egg?

- After 4 years – back to the beginning
  - Applying theory/methodology more fully
  - Applying at the beginning of integrations
  - Applying more comprehensively across the warehouse
- Major re-thinking of all Galaxy processes
  - Interacting with customers
  - Writing specifications
  - Obtain source files
  - Mapping source-to-target
  - Implementing ETL
  - Building physical tables
  - Taking DQ measures
- DQ still requires championing
- New problem: How to analyze and respond to findings from the data gathered through new process.

**INGENIX.**