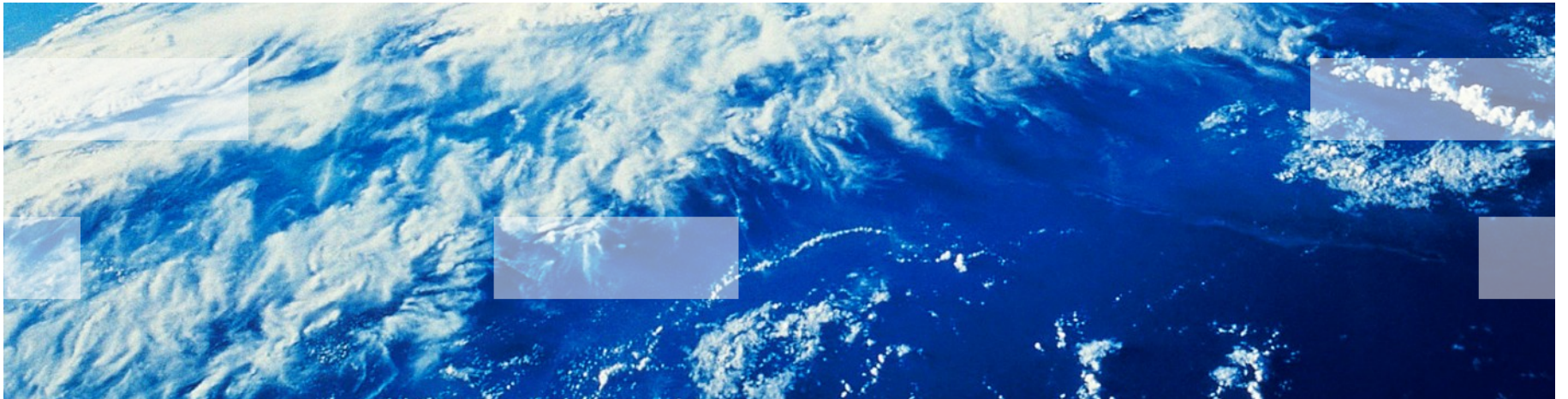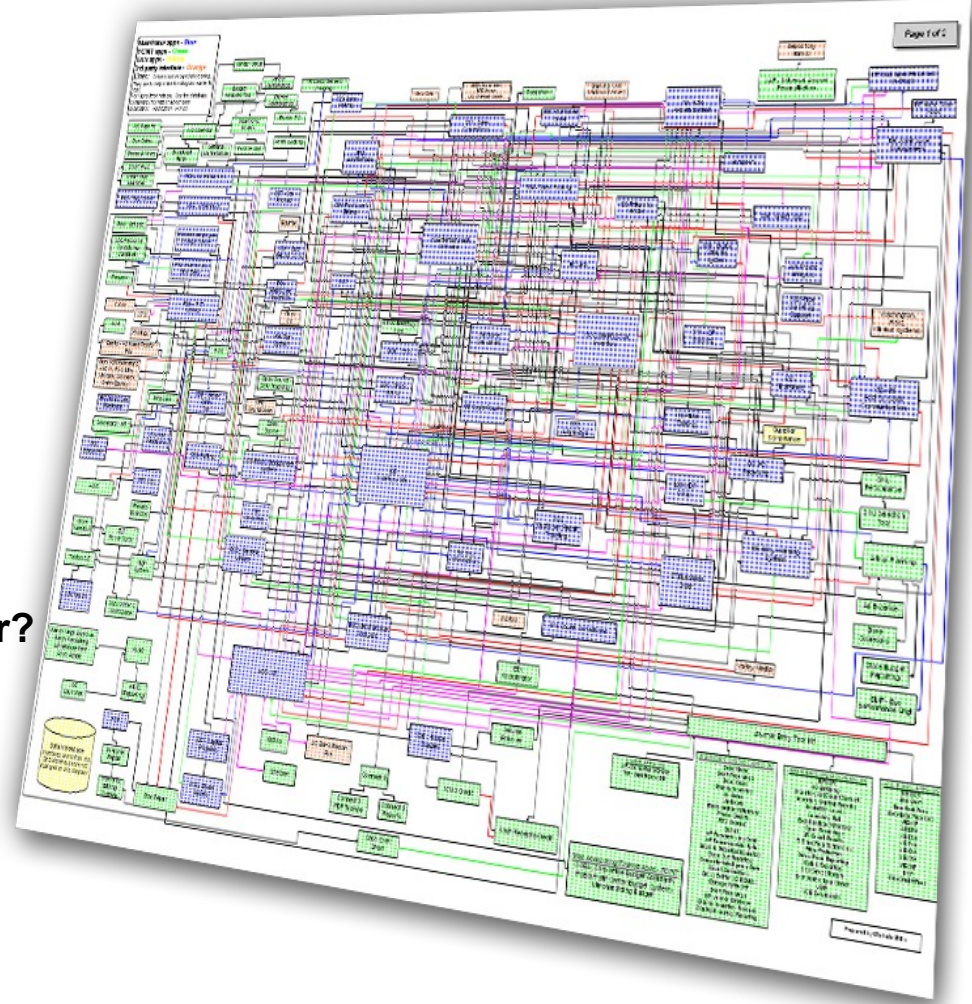IBM yazılım zirvesi '09 style

Tevfik Baskaya – Vodafone & Okay Akyüz - OSİS
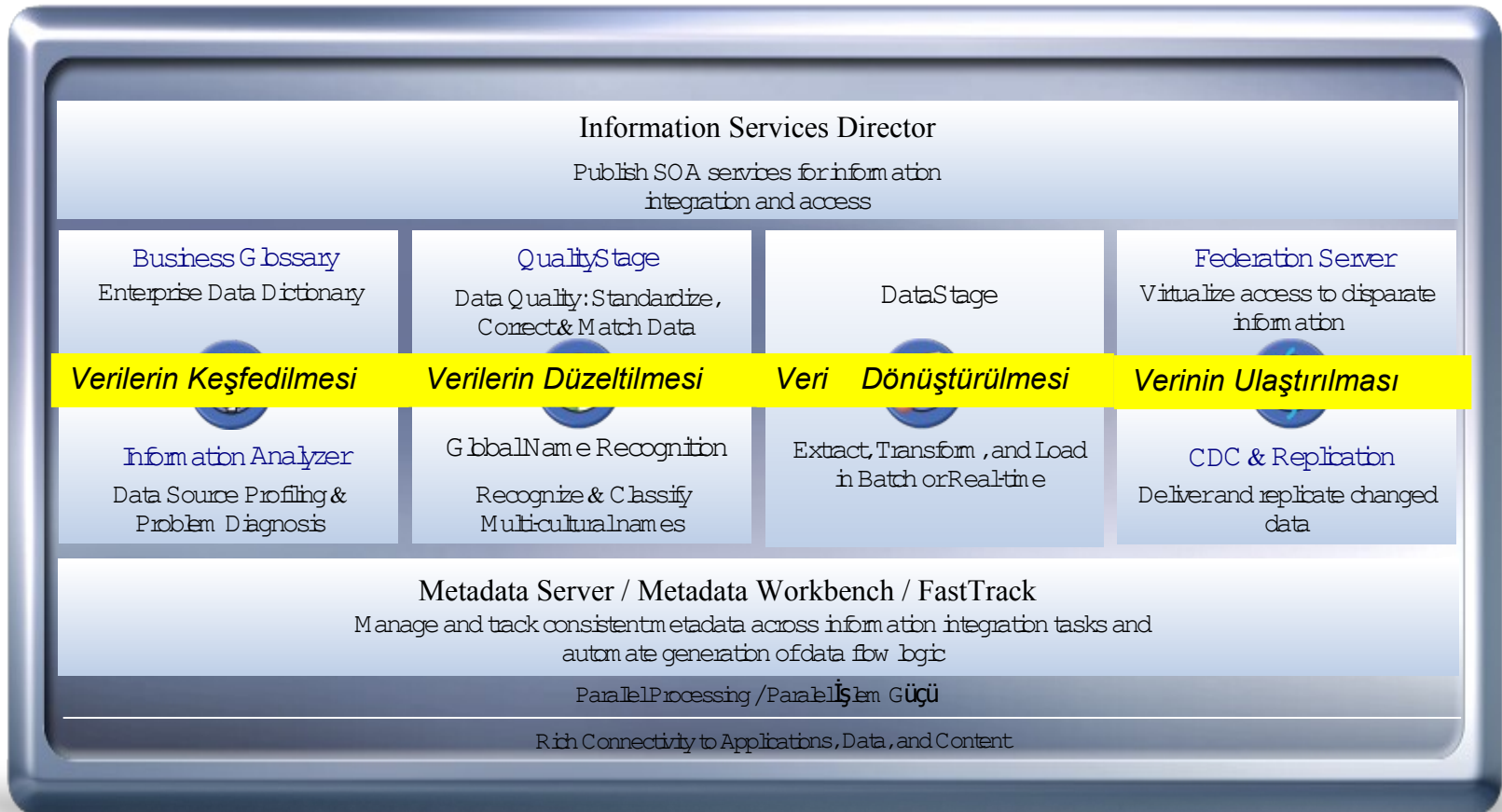
22/10/2009

# Bilgi...

- **Bilgilerim nerede?**

- **İhtiyacım olduğunda nasıl erişirim?**

- **Bilgilerimin içeriğinde ne var?**

- **Güvenilir mi?**

- **Kalitesini nasıl arttırabilirim?**

- **İstediğim şekilde nasıl erişirim?**

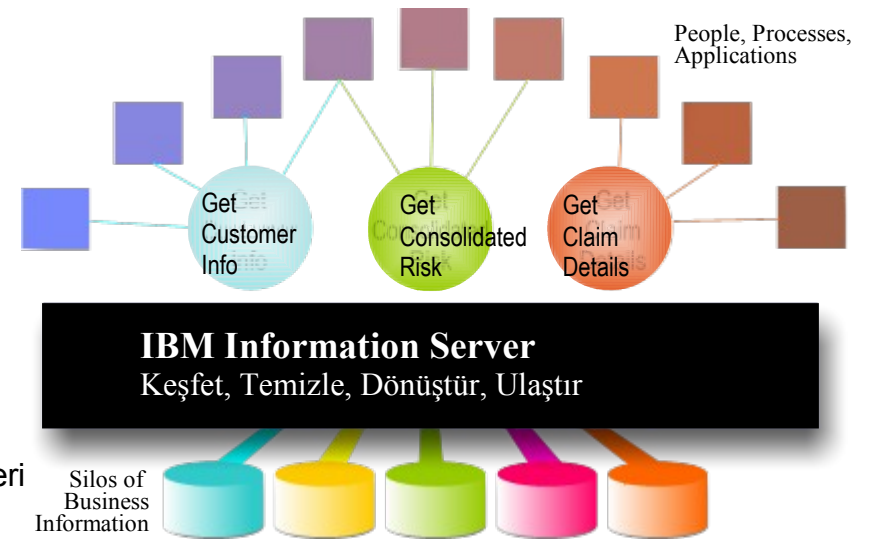- **Bilgim nerden geliyor, nereye gidiyor?**

- **Nasıl kontrol edebilirim?**

# InfoSphere Information Server – Ürünleri

**Information Services Director**

Publish SOA services for information
integration and access

| Business Glossary | QualityStage | DataStage | Federation Server |
|---|---|---|---|
| Enterprise Data Dictionary | Data Quality: Standardize, Correct & Match Data | | Virtualize access to disparate information |
| **Verilerin Keşfedilmesi** | **Verilerin Düzeltilmesi** | **Veri Dönüştürülmesi** | **Verinin Ulaştırılması** |
| Information Analyzer | Global Name Recognition | Extract, Transform, and Load in Batch or Real-time | CDC & Replication |
| Data Source Profiling & Problem Diagnosis | Recognize & Classify Multi-cultural names | | Deliver and replicate changed data |

**Metadata Server / Metadata Workbench / FastTrack**

Manage and track consistent metadata across information integration tasks and
automate generation of data flow logic

Parallel Processing / Paralel İşlem Gücü

Rich Connectivity to Applications, Data, and Content

33

# IBM Information Server Bileşenleri

- ## Information Analyzer
  - Veri Analizi, raporlama

- ## DataStage
  - Verilerin ETL süreçlerini yönetir

- ## QualityStage
  - Veri standartlaştırması ve temizliği

- ## Bussiness Glossary
  - Veri Sahipliği, Metadata Yönetimi

- ## Metadata Workbench
  - Metadata Arayüzü, Etki Analizleri, Veri Akışı Analizleri

- ## Fast Track
  - Hızlı ETL geliştirme ve proje yönetimi

- ## Change Data Capture
  - Değişen veriyi yakalama, realtime datawarehouse

- ## Information Services Director
  - Web Servis ile kurumsal veri katmanı oluşturulması

People, Processes, Applications

Get Customer Info

Get Consolidated Risk

Get Claim Details

**IBM Information Server**
Keşfet, Temizle, Dönüştür, Ulaştır

Silos of Business Information

44

# IBM Information Server – Business Glossary

- Web tabanlı iş ana verisinin paylaşımı, yönetimi ve yazılması

- IT'nin çalışmalarını iş hedefleri ile eşleştirir
- IT değerlerine iş içeriğinin eklenmesini sağlar

- Sahiplik ve sorumluluk tanımlamalarını yapar

**Subject Matter Experts**

**Business Users**

**Understand**

**Business Glossary**
İş terimleri kütüphanesini yaratır ve yönetir. İş terimlerinin ilişkilerini tanımlar, fiziksel tasarım ile iş terimlerini ilişkilendirir

Database = DB2
Schema = NAACCT
Table = DLYTRANS
Column = ACCT_NO
data type = char(11)

Müş. Hesap No On haneli hesap numarası. ACCNO ve HESNO diye de adlandırılır L-FIIIIVVVV formatında yazılır

eporting                                          Help | Log O

Overview | Manage Custom Attributes × | Manage Business Terms ×

Overview

**Welcome to WebSphere Business Glossary**
International Business Machines (NYSE:IBM) is the leader in enterprise data integration. Customers and partners worldwide use its information integration tools to confidently transform data into accurate, reliable and complete business information to improve operational performance and decision-making across every critical business dimension. Our comprehensive end-to-end solutions provide on demand data integration complemented by our professional services, industry expertise, and methodologies.

IBM.

USA
Finance
Manufacturing
Customer Relationship

Europe
Finance
Manufacturing
Customer Relationship

Data Model
Top-level category, which groups all the elements by type

Asia
Finance
Manufacturing
Customer Relationship

See All Categories

**İş Odaklı Görünüm**

## Business Glossary – Neden İhtiyaç?

- Kurumlardaki uzmanlıklar veriyi ve problemleri anlamakta, sonuçları yorumlamakta kritik rol oynarlar
- Yetersiz uzmanlık ve bilgi birikimi düşük veri kalitesinin başlıca nedenidir, veri kullanılmaz hale dönüşür
- Genellikle çalışanların hafızalarındadır – nadiren dokümante edilir
- Kurum içerisinde dağılmış haldedir
- İşten ayrılmalarda ve proje değişimlerinde kaybolur
- Dokümante edilmemişse, bilgi ve yorumlama bozulmakta ve zamanla bulanıklaşmaktadır

Developer

Business Analyst

End Users

Software Architect

Data Architect

Data Admin

IT Admin

# Business Glossary

**Yararları: Fikirler, düşünceler ve bilgi birikimlerinin paylaşımı ile iş birliğini yaratır**

**Muhasebe Bölüm Kodu**

SAHIBI: Muhasebe Bölümü
FORMAT: X(7)
TANIM: Hesabın bağlı olduğu organizasyon bölümünü tanımlayan yedi haneli numara

**Standart Tanımı Yap**

**Herhangi bir nesneye not ekle**

Alt bölümü tanımlayan Bölüm Kodunun son 2 iki hanesinin genelde boş olduğunu farkettim

# IBM Information Server – Business Glossary

## Bilgi Varlıklarınızın Tamamına Hakim Olun



**Business Glossary**

- İşinizi tanımlayan terimler nelerdir?

- Ne anlam ifade ederler?

- Bu terimler birbirleri ile bağlantıları nelerdir?

- Bu terimlerin iş sahipleri kimler?

**Birleşik Değer**

- Bu sistemdeki veri ne anlama geliyor?
- Bu sistemdeki iş terimlerinin kuralları nedir?

**Information Analyzer**

- Bu sistemdeki verinin yapısı nasıl?

- Verinin biçimi nedir?

- Veri kalitesi nasıl?

- Bu iki sistem nasıl bağlantılılar?

# IBM Information Server – Business Glossary

- Business Glossary olmayan firmalarda bütün iş sözlüğü bir excel listesinde tutulmaktadır.

- Business Analistlerin ve bütün kullanıcıların iş sözlüğüne erişmeleri saplanıyor

- Herhangi bir metrik'e ihtiyaçımız olursa bunu ortak common metadata sayesinde iş sözlüğün'de bakıyoruz

- Herhangi bir metrik'i arayıp buluruz ve bu metrik ile ilgili tüm detaylar business – glossary'de mevcuttur.

  - Short description
  - Long description
  - Steward
  - Related It assets
  - Etc.

# Information Analyzer

- Varolan sisteminizi derinlemesine analiz eder

    - Veritabanları ve dosya bazlı kaynakların içerik, kalite ve yapısal olarak veri odaklı analizini yapar
    - Kolonların, kolonlar arasındaki ve veritabanları arasındaki ilişkilerin detaylı profilini çıkartır

- Veri kalitesini, önceki raporlarla karşılaştırarak sürekli olarak gözlemlenebilmesini sağlar

- Bilginin kurum içinde nerede ve nasıl yönetildiğinin ana veri bilgisini oluşturur

    - Kurum içindeki uzmanlık bilgisinin de eklenmesi ile, yeni projelerin daha hızlı yapılmasına ve doğru bilgi ile yapılmasını sağlar

**Subject Matter Experts**

**Data Analysts**

**Information Analyzer**

**Understand**

**Kaynak veri yapılarını inceler, veri katilesini ve bütünlüğünü gözlemler**



**Fiziksel Görünüm**

# IBM Information Server – Information Analyzer

- ## Information Analyzer Fonksiyonel Bakış

  - The Project Paradigm
    - Project Artifacts & Elements
  - Managing Data
    - Platform Subject Administration
    - Metadata Import
  - Profiling & Review
    - Column Analysis
    - Primary Key Analysis
    - Foreign Key Analysis
    - Redundancy Analysis
    - Cross Domain Analysis
    - Baseline Analysis
    - Sharing Results
  - Developing Data Rules & Metrics
    - Data Rules
    - Benchmarks
    - Metrics

11

# Information Analyzer : Kolon Analizi

# Information Analyzer : Kolon Analizi

M Corporation

# Information Analyzer : Kolon Analizi

# Information Analyzer : Kolon Analizi

# Information Analyzer : PK, FK Analizi

# Tablolar Arası/Sistemler Arası Seviyesinde Bilgilendirme



**İlişkisel Bütünlük**

- Foreign Keys
  - Tekil veya Çoğul Alan
  - Referential Integrity

- Cross-Domain
  - Benzerlik
  - Redundancy

18

# QualityStage

- Bütünleşik Veri kalitesi süreci sunar
  - Temiz, standart ve tekil veriye sahip olmanızı sağlar
  - Gerçeğin tek bir versiyonunu oluşturur
  - Global Adres doğrulama

- Veri kalitesi ve eşleştirme algoritmaları için görsel araç
  - Quality Stage ile tam entegre (tek motor, tek veri modeli, tek arayüz)
  - Eşleştirme kurallarını kolay ve kesin olarak ayarlama

- Veri kalitesi mantığını ETL içine kolayca ekler

**Subject Matter Experts**

**Data Analysts**

**Cleanse**

**QualityStage™**
**Alanları standartlaştırır ve düzeltir, değişik kaynaklardaki kayıtları eşleştirerek verinin tek bir görüntüsünü oluşturur**

**Visual Match Rule Design**

# QualityStage

## Genel Veri Problemleri

· **Verinin standartlaşma eksikliği**
  – Değişik sistemlerde değişik şekillerde ve yapılarda durmakta

· **Alanlardaki farklı veriler**
  – Veri veritabanında yanlış kolonlarda durmakta

· **Bilgi, serbest metin içinde kaybolmakta**

· **Veri miyopluğu**
  – Farklı kaynaklarda duran, ilişkilendirilemeyen ilişkili veriler

· **Veri tekrarı**
  – Standartları farklı kopya bilgiler

```
Kate A. Roberts    416 Columbus Ave #2, Boston, Mass 02116

Catherine Roberts Four sixteen Columbus APT2, Boston, MA 02116

Mrs. K. Roberts    416 Columbus Suite #2, Suffolk County 02116
```

```
Name            Tax ID          Telephone

J Smith DBA Lime Cons.  228-02-1975    6173380300
Williams & Co. C/O Bill  025-37-1888    415-392-2000
1st Natl Provident      34-2671434     3380321
HP 15 State St.    508-466-1200  Orlando
```

```
WING ASSY DRILL 4 HOLE USE 5J868A HEXBOLT 1/4 INCH

WING ASSEMBY, USE 5J868-A HEX BOLT .25" - DRILL FOUR HOLES

USE 4 5J868A BOLTS (HEX .25) - DRILL HOLES FOR EA ON WING ASSEM

RUDER, TAP 6 WHOLES, SECURE W/KL2301 RIVETS (10 CM)
```
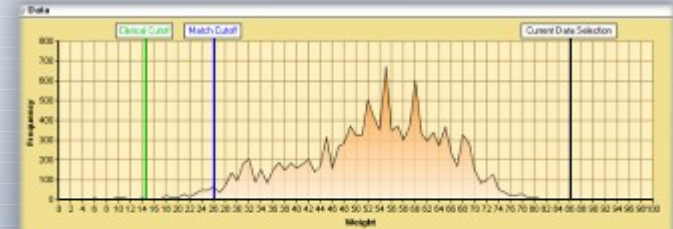
```
19-84-103    RS232 Cable 6' M-F CandS

CS-89641     6 ft. Cable Male-F, RS232 #87951

C&SUCH6      Male/Female 25 PIN 6 Foot Cable
```

```
90328574   IBM              187 N.Pk. Str. Salem NH 01456
90328575   I.B.M. Inc.          187 N.Pk. St. Salem NH 01456
90238495   Int. Bus. Machines   187 No. Park St Salem NH 04156
90233479   International Bus. M. 187  Park Ave Salem NH 04156
90233489   Inter-Nation Consults 15 Main Street Andover MA 02341
90345672   I.B. Manufacturing    Park Blvd. Bostno MA  04106
```

## İşinizde Tutarlı, Doğru ve Birleştirilmiş Veriyi Nasıl Elde Edeceksiniz?

**Müşteriler**

**Ürünler / Malzemeler**

**İşlemler**

**Üreticiler / Tedarikçiler**

**1. Araştırma (Investigatio**

**QualityStage Süreçleri**

**Hedef**

Veritabanında Birleştirilmiş Görüntüler

21

# InfoSphere Information Server – Ürünleri



**Information Services Director**

Publish SOA services for information
integration and access

**Business Glossary**
Enterprise Data Dictionary

**QualityStage**
Data Quality: Standardize,
Correct & Match Data

**DataStage**

**Federation Server**
Virtualize access to disparate
information

**Information Analyzer**
Data Source Profiling &
Problem Diagnosis

**Global Name Recognition**
Recognize & Classify
Multi-cultural names

Extract, Transform, and Load
in Batch or Real-time

**CDC & Replication**
Deliver and replicate changed
data

Metadata Server / Metadata Workbench / FastTrack
Manage and track consistent metadata across information integration tasks and
automate generation of data flow logic

Parallel Processing / Paralel İşlem Gücü

Rich Connectivity to Applications, Data, and Content

# IBM yazılım zirvesi '09 style

# Vision

With this project, we focused on two main capabilities

1. Data Profiling
2. Automating Data Quality Check with flexible and reportable( open to users) applications.

# Data Profiling

- Telecommunication Data
    - Huge Data size and growing fast
    - New data and concepts dynamic business
    - New coming applications in portfolio.

- During journey of information so many applications and systems have integration points. <span style="color:red">So high risk for inconsistent information</span>.

- By Using Data Profiling tool
    - Identify Anomalies
    - Identify Inconsistencies
    - Taking Corrective Actions – Planning DQ approach.

# Why Data Quality Tool

1. Operational Efficiency

1. Maintainability & Manageability

1. Low Impact to Source Systems

1. Performance & Scalability

1. Open Environment – Reports available to business users

# Selection Criterias

- Vodafone is again pioneer; writing Data Quality Jobs that checks data in DWH is somehow new in Turkey.

- So we needed a strong tool that is supported by strong Vendor.

- Product must be well-supported against problems.

- We needed a strong solution partner.

- Solution partner should have deep knowledge both in technology and customer's business.
  - I

# Selection Criterias

· Performance

· Easy to develop/modify applications via a GUI

· Interface/Connectivity to other systems

· Integration with other DWH tools

· Vendor/Support

· Solution Integrator
  - I

# Project Actions

- Analysis for data to be checked – Identifying DQ Controls
- Thresholds for each data and controls
- Identification corresponding actions (alarms, mail)
- Identification of reporting requirements
- Development for Phase 1
- Pilot run
- Production (Not Yet)

Current status is pilot run completed in IT for Phase 1.

- I

# General Requirements

- Flexible

  - Parametric threshold definition for each control

  - Parametric action definition

  - Easy to change thresholds and actions

- Manageable & Reportable

  - Results are stored in Relational Database

  - Results are accessed via reporting interface, open to all users

  - Results are stored historically to show trends , collaborative

- All jobs should use the same framework.

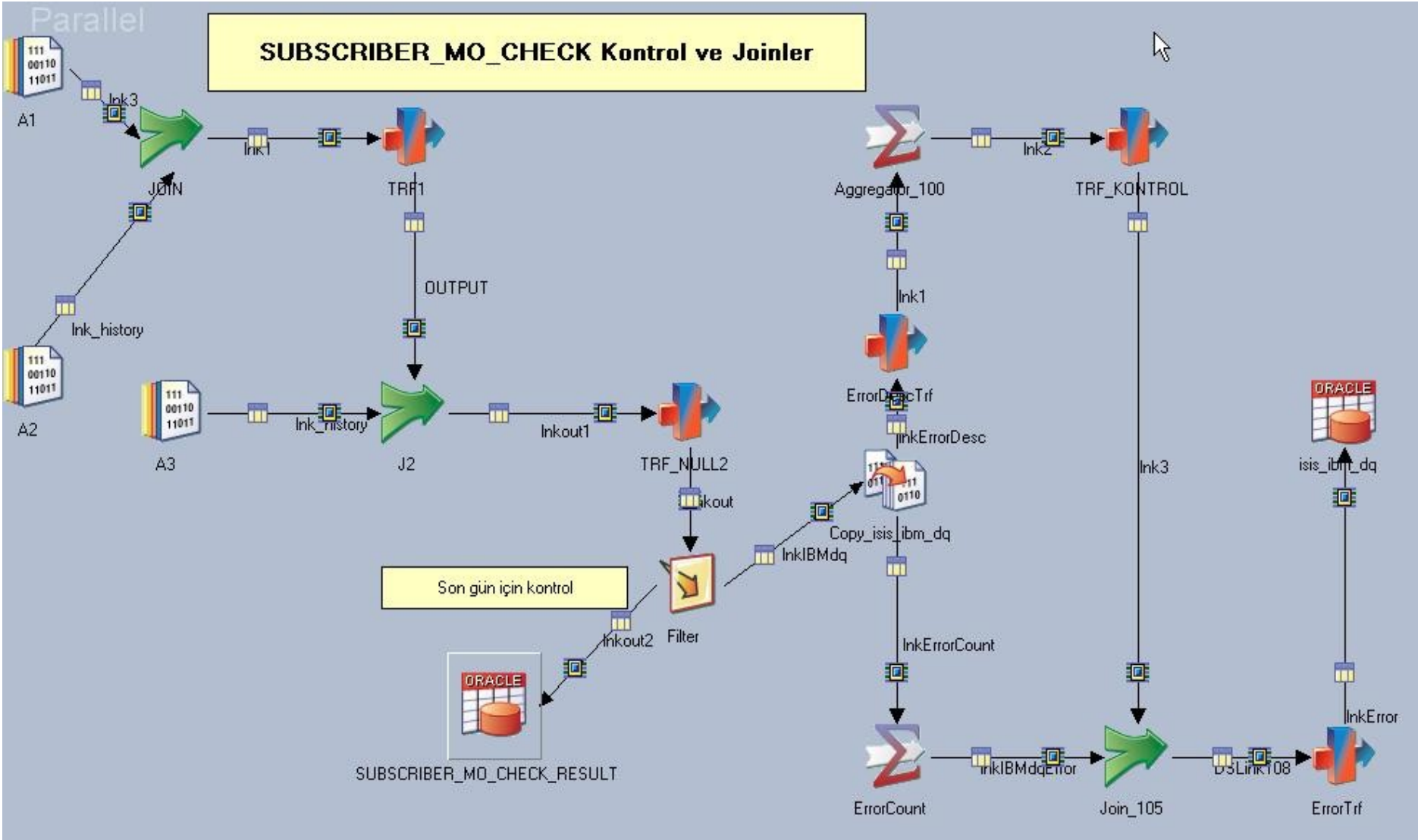- Opening results & trends to users produce auditable framework.

  - I

# Sample Case

| Day | Snapshot | Correct |
|-----|----------|---------|
| 9 Oct | 5 Oct | 5 Oct |
| 8 Oct | 5 Oct | 5 Oct |
| 7 Oct | 5 Oct | 5 Oct |
| 6 Oct | 5 Oct | 5 Oct |
| 5 Oct | 5 Oct | 5 Oct |

| Day | Snapshot | Correct |
|-----|----------|---------|
| 10 Oct | 10 Oct | 10 Oct |
| 9 Oct | 5 Oct | 8 Oct |
| 8 Oct | 5 Oct | 8 Oct |
| 7 Oct | 5 Oct | 5 Oct |
| 6 Oct | 5 Oct | 5 Oct |
| 5 Oct | 5 Oct | 5 Oct |

# Örnek uygulama



SUBSCRIBER_MO_CHECK Kontrol ve Joinler

*Bu sunum 22 Ekim 2009 tarihinde İstanbul  Swissotel the Bosphorus'da yapılan Yazılım Zirvesi 2009 için hazırlanmıştır.*

*http://www.ibm.com/software/tr*