

Технологии обработки информации

Лекция 8-9

Интеграция информационных ресурсов

Содержание

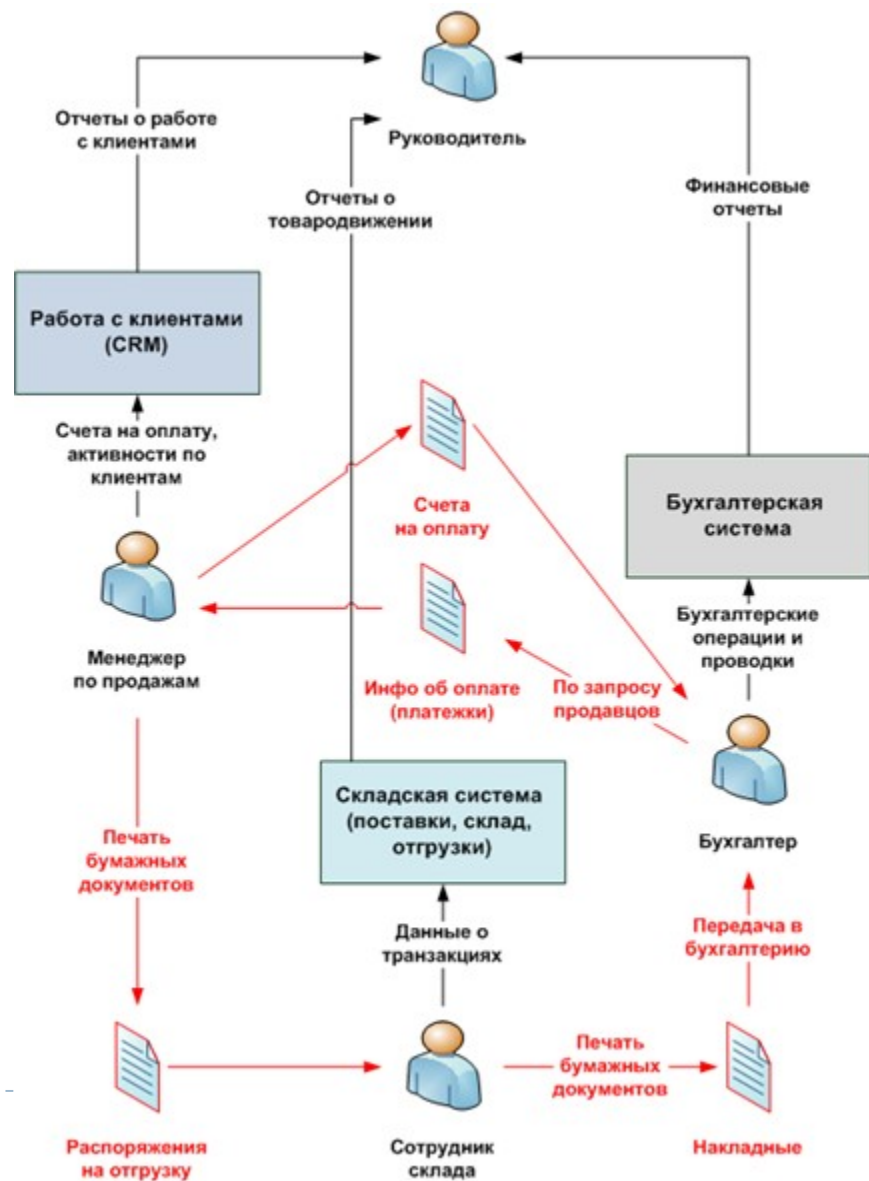
- ▶ Проблема интеграции данных
- ▶ Классификации методов интеграции
- ▶ Что такое SQL Server 2005 Integration Services
- ▶ Планирование ETL проекта

Интеграция



Что, куда, зачем?

Системная интеграция: проблема



Самые востребованные ИТ-профессии-2012 в мире

1. **Разработчик бизнес-архитектуры** - занимается моделированием совместимости главных стратегий компании с технологиями
2. **Специалист по информации** - работает с большим потоком неструктурированной или полуструктурированной информации, которая получена из самых разных источников, включая веб-страницы, журналы учета заданий и т.д. Для такой работы необходимы сотрудники с разнообразными умениями от подготовки данных для анализа до обработки статистики

Computerworld.com

Интеграция данных

- ▶ включает объединение данных, находящихся в различных источниках и предоставление данных пользователям в унифицированном виде.
- ▶ Уровни интеграции:
 - ▶ **физический** - конверсия данных из различных источников в требуемый единый формат их физического представления
 - ▶ **логический** – создание единой глобальной схемы данных
 - ▶ **семантический** - единого представления данных с учетом их семантических свойств в контексте единой онтологии предметной области

Проблема интеграции данных

ETL процессы (Extraction, Transformation, Load)

60 -80% времени

- Извлечение и очистка данных
- Трансформации данных
- Загрузка данных

Типы несоответствия схем данных

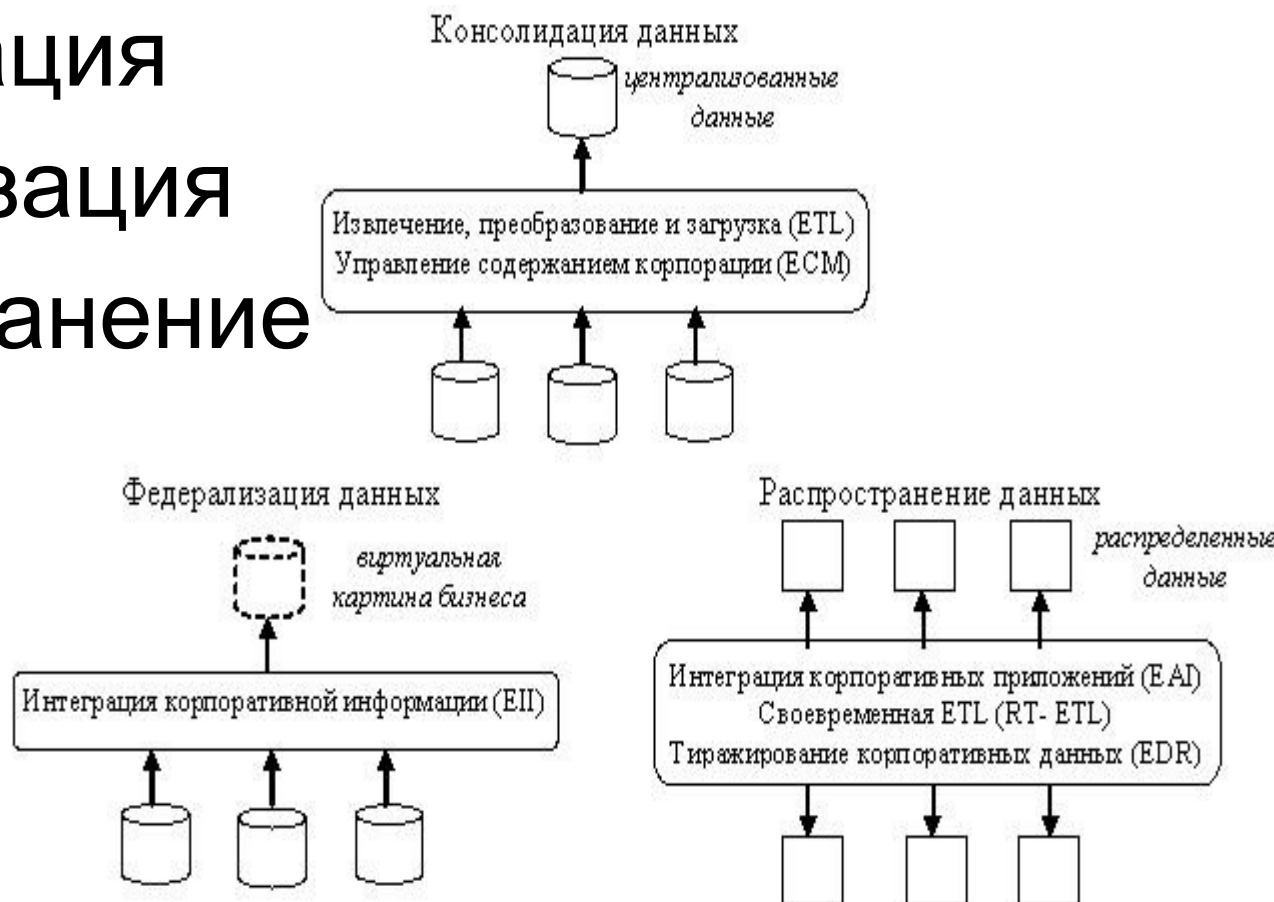
- ▶ Конфликты неоднородности (используются различные модели данных для различных источников)
- ▶ Конфликты именования (в различных схемах используется различная терминология, что приводит к омонимии и синонимии в именованиях)
- ▶ Семантические конфликты (выбраны различные уровни абстракции для моделирования подобных сущностей реального мира)
- ▶ Структурные конфликты (одни и те же сущности представляются в разных источниках разными структурами данных)

Типы несоответствия собственно данных

- ▶ Различие формата данных
- ▶ Различие в представлении значений
- ▶ Потеря актуальности данных одним из источников
- ▶ Наличие ошибок операторского ввода (или ошибок распознавания бланков) в отдельных источниках данных
- ▶ Намеренное внесение искажений с целью затруднить идентификацию сущностей

Три метода интеграции данных

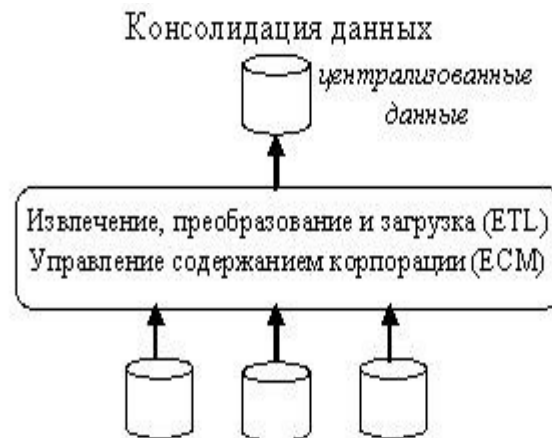
- ▶ Консолидация
- ▶ Федерализация
- ▶ Распространение



- ▶ SOA
- ▶ семантическая интеграция

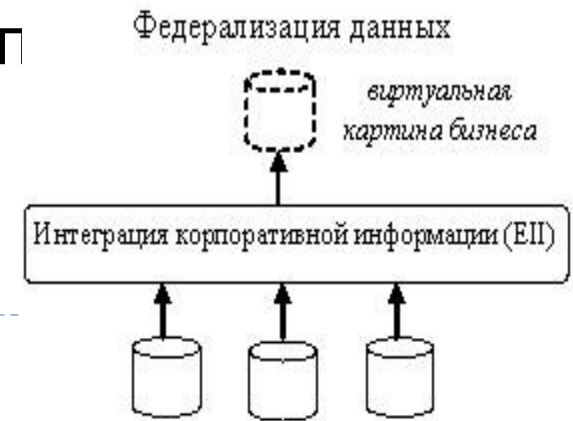
Консолидация данных

- ▶ Данные собираются из нескольких первичных систем и интегрируются в одно постоянное место хранения. Такое место хранения может быть использовано для подготовки отчетности и проведения анализа, как в случае хранилища данных, или как источник данных для других приложений.



Федерализация данных

- Обеспечивает единую виртуальную картину нескольких первичных источников данных. Для получения сведений о некотором процессе, обрабатываемом в нескольких оперативных приложениях, процессор федерализации данных извлекает данные из соответствующих первичных складов данных, интегрирует их таким образом, чтобы они отвечали виртуальной картине и требованиям запроса, и отправляет результаты бизнес-приложению, от которого п



Распространение данных

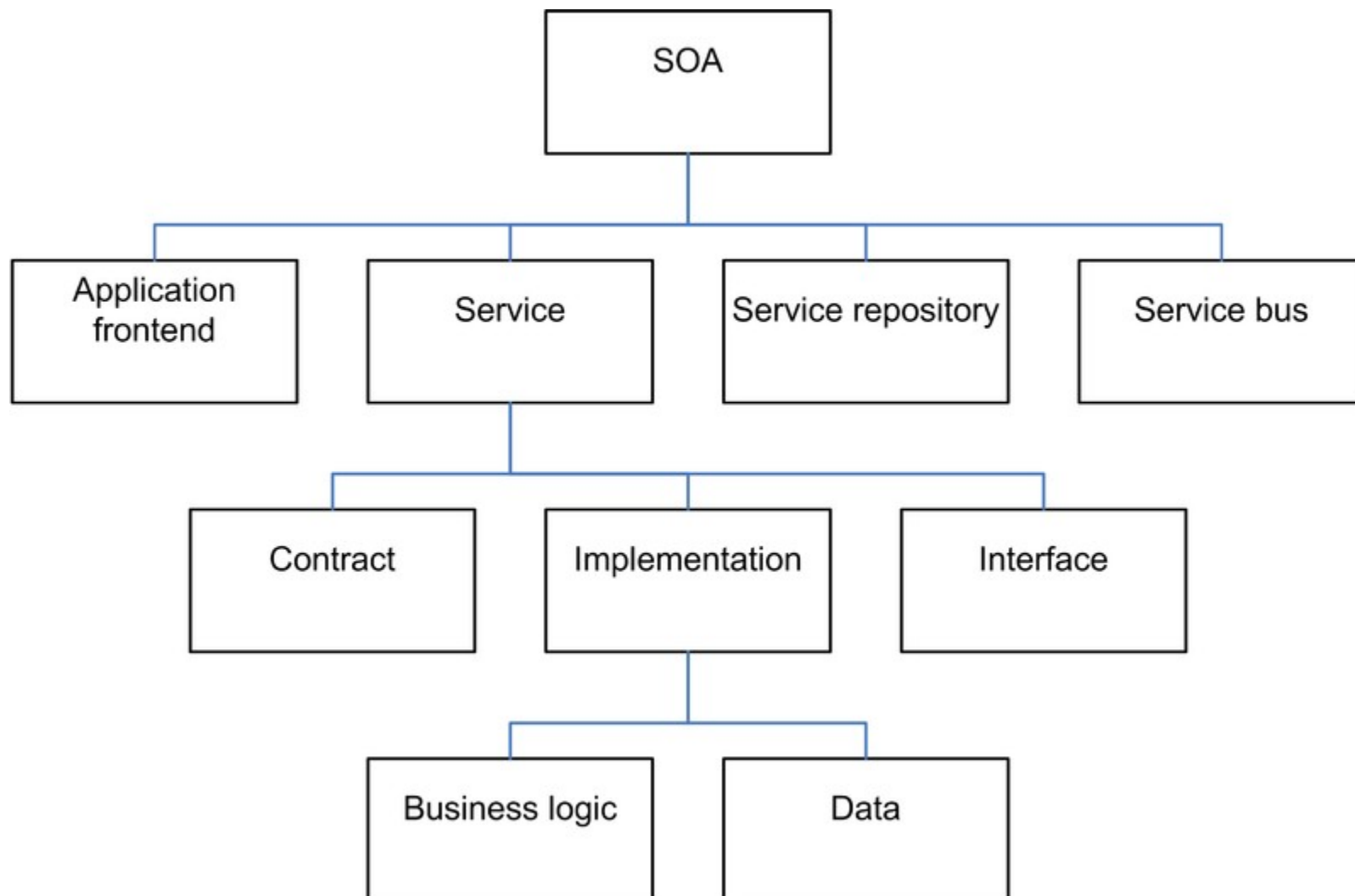
- ▶ Подразумевает их копирование из одного места в другое. Этот подход обычно используется для операций реального времени и базируется на механизмах "проталкивания", т. е. является событийно управляемым.



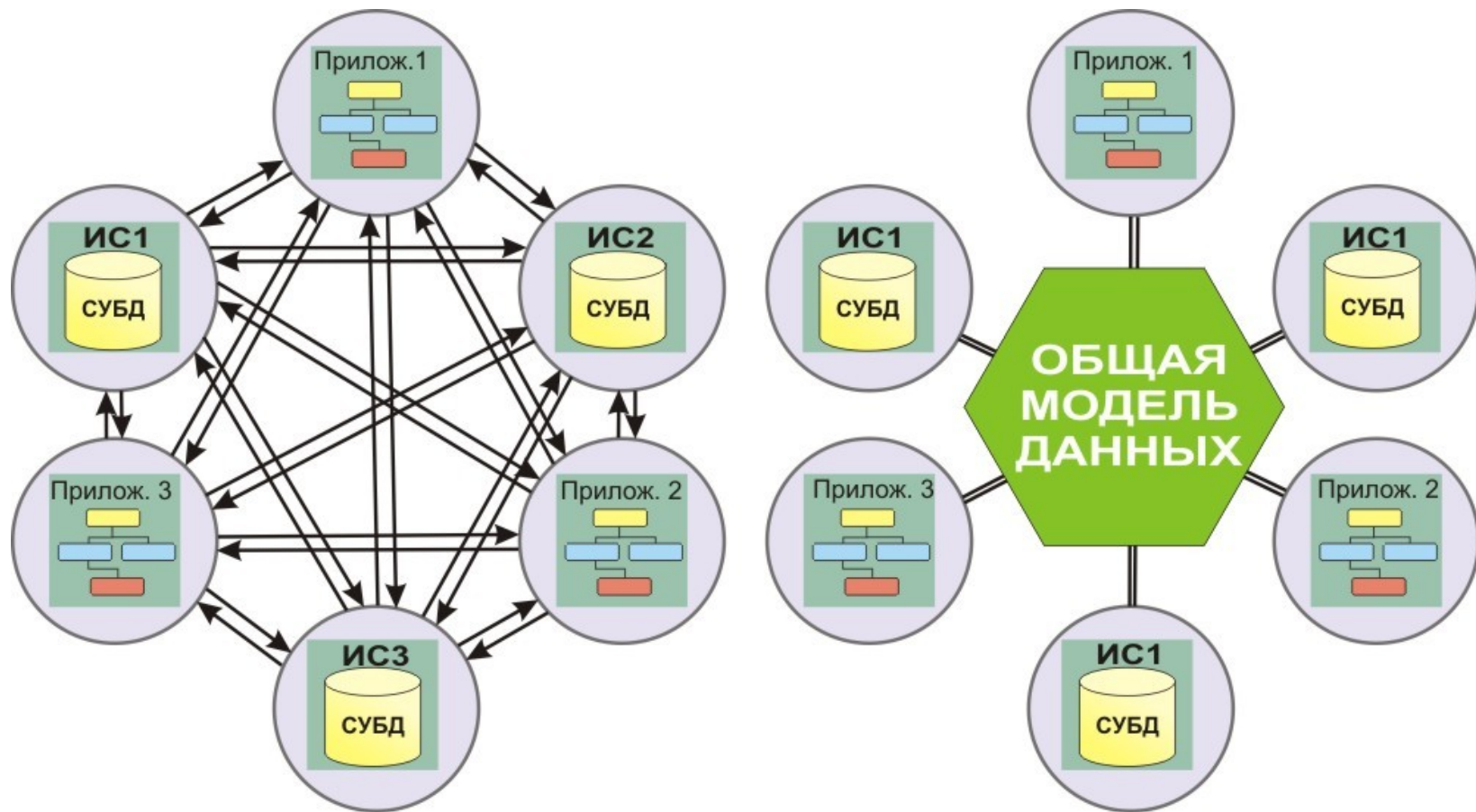
Сервисный подход – SOA (1)

- ▶ Service Oriented Architecture
- ▶ модульный подход к разработке программного обеспечения, основанный на использовании распределённых, слабо связанных (англ. *loose coupling*) заменяемых компонентов, оснащённых стандартизированными интерфейсами для взаимодействия по стандартизированным протоколам
- ▶ Данные остаются у владельцев и даже местонахождение данных неизвестно
- ▶ При запросе происходит обращение к определённым сервисам, которые связаны с источниками, где находится информация и ее конкретный адрес

Сервисный подход – SOA (2)



Интеграция на основе метамодели (семантическая интеграция)



Классификация методов интеграции данных по Клаусу Диттриху



<http://www.osp.ru/os/2009/10/11170978/>

Задачи при интеграции данных

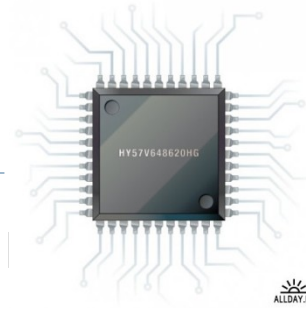
- ▶ Технологические
- ▶ Организационные
- ▶ Экономические



ALLDAY.RU



Технологические задачи

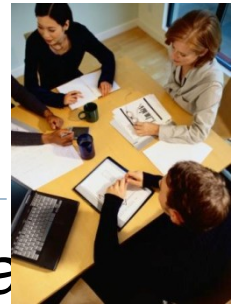


- ▶ Гетерогенные источники данных с различными форматами
- ▶ Структурированные, полуструктурированные и неструктурированные данные
- ▶ Данные поступают в разное время
- ▶ Очень большие объемы данных
- ▶ Качество данных (пропуски, нет смысла, ошибки)
- ▶ Придание смысла данным при слиянии их из разных форматов при неполноте данных в отдельных источниках
- ▶ Преобразование данных в унифицированный формат, пригодный для бизнес-анализа

Технологические требования

- ▶ Загрузка данных в наибо́льшее время (нет возможности «ночного» периода, 7 x 24 часа On-Line)
- ▶ Потребность загрузки данных в несколько приемников практически одновременно
- ▶ Постоянная доступность данных с минимальными задержками в актуальности данных
- ▶ Разнообразие источников данных (OLTP, OLAP, веб-сервисы, неструктурированные данные, унаследованные системы)
- ▶ Разнообразие приемников данных (порталы, персонализированные отчеты, PDA, мобильные телефоны)

Организационные задачи



- ▶ Получение серьезной поддержки руководства компании команде по проекту интеграции данных, настоять на координации и компромиссах по выбору форматов данных и бизнес-процессов получения данных в подразделениях компании
- ▶ Определиться с единообразными технологиями для разного круга задач, так как многие подразделения используют совершенно разные системы и способы. Люди консервативны в своих привычках, не любят переучиваться. До 60% времени при получении и интеграции данных — ручной процесс

Экономические задачи



Интеграция данных – дорогостоящий процесс.

Факторы, увеличивающие стоимость проекта:

- ▶ Административные преграды, недостаток координации, недостаточная поддержка руководства
- ▶ Недостаточная функциональность имеющихся средств для ETL процессов, необходимость разработки нового ETL кода

Интеграционные платформы

- ▶ Microsoft BizTalk Server
- ▶ Microsoft SQL Server
- ▶ Oracle SOA Suite
- ▶ IBM WebSphere

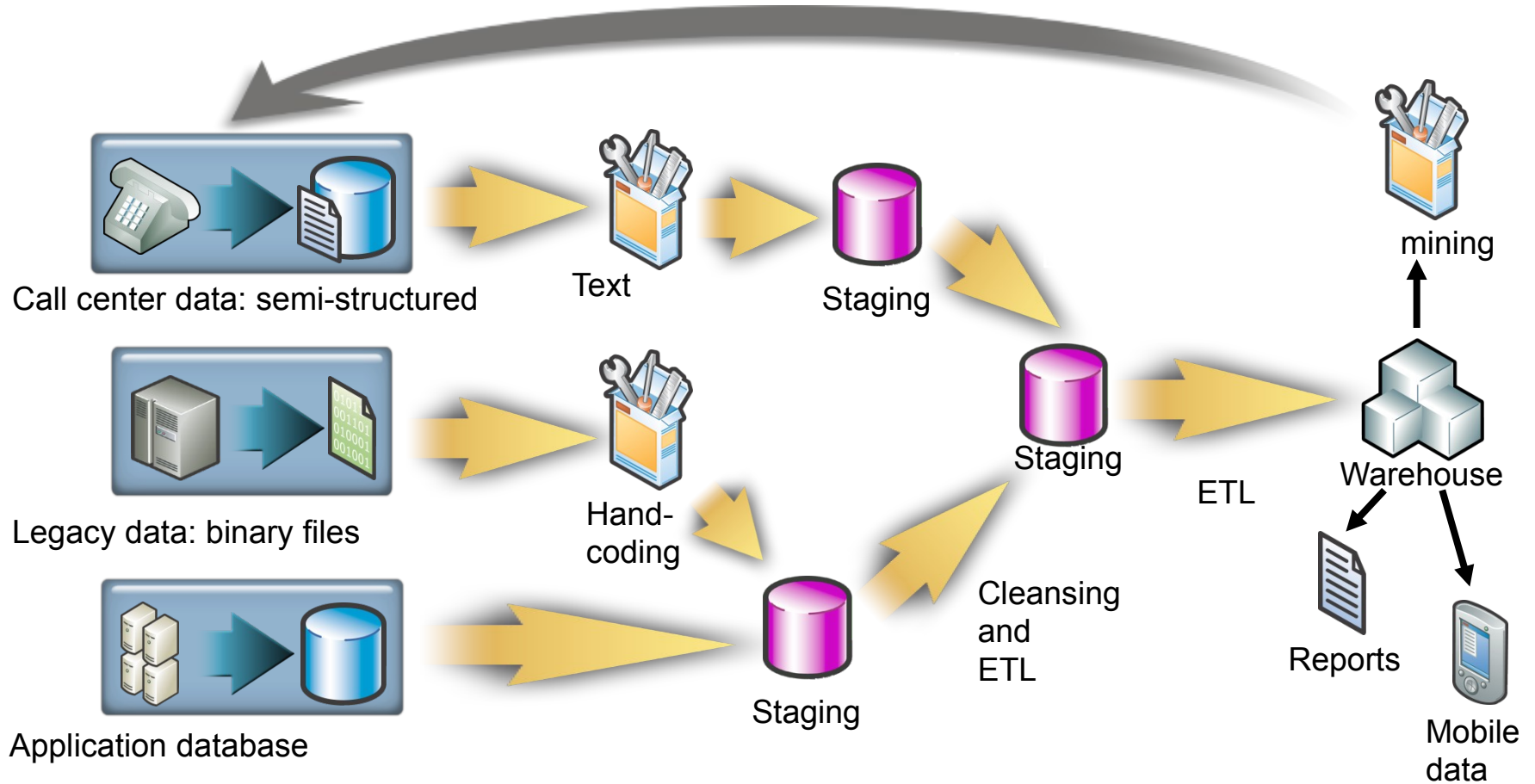


SQL Server 2008 Integration Services

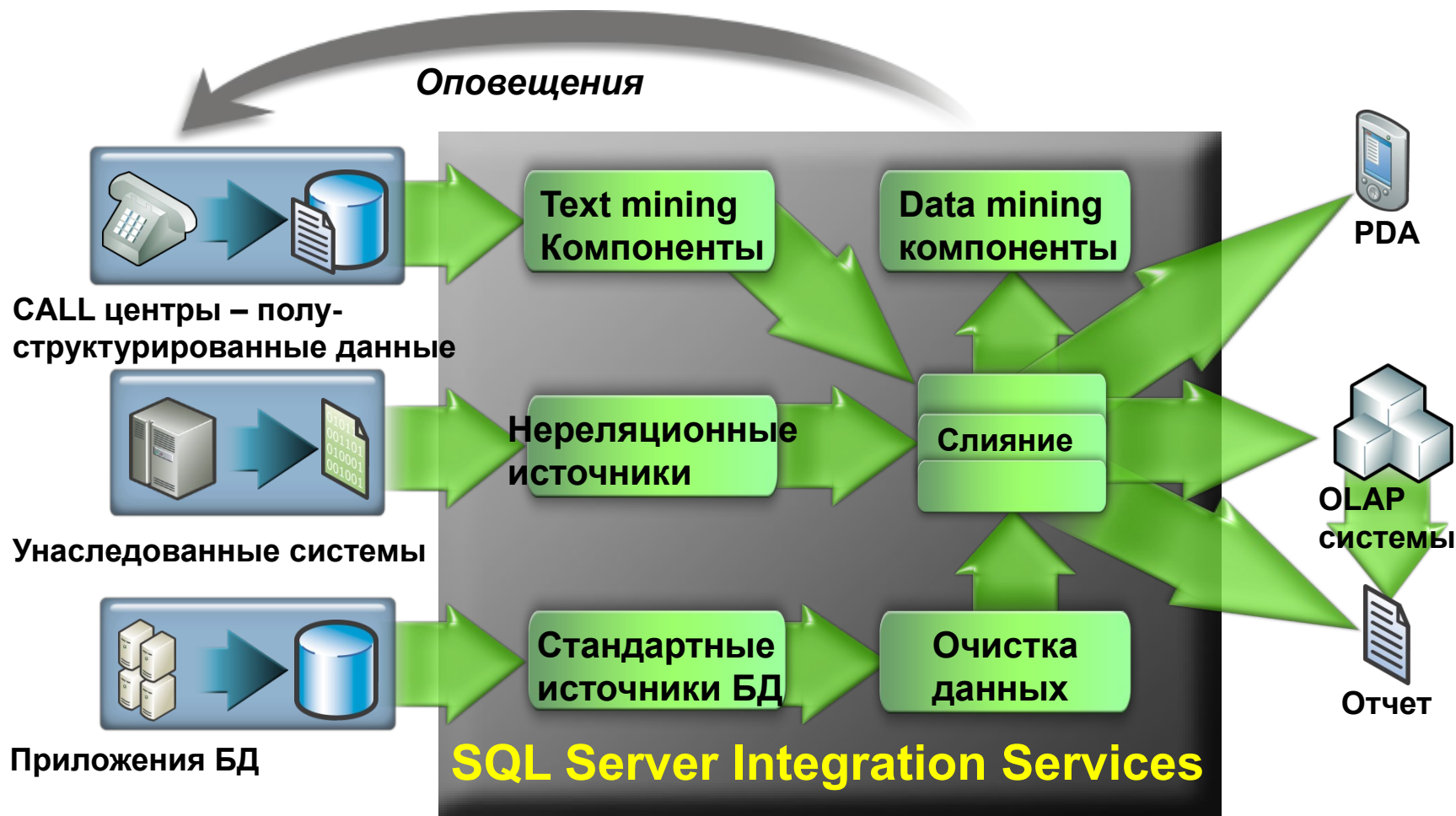
Службы Integration Services - платформа для построения высокопроизводительных решений интеграции данных и решений потока операций, включая операции извлечения, преобразования и загрузки (ETL) для хранилищ данных.

- Графические инструменты
- Мастера для построения и отладки пакетов
- Источники данных для извлечения данных
- Источники назначения для загрузки данных
- Преобразования для очистки, статистической обработки, слияния и копирования данных
- Задачи для выполнения функций потока операций
- Служба управления и администрирования пакетов
- API-интерфейсы для программирования объектной модели

Do Integration Services



Integration Services 2008



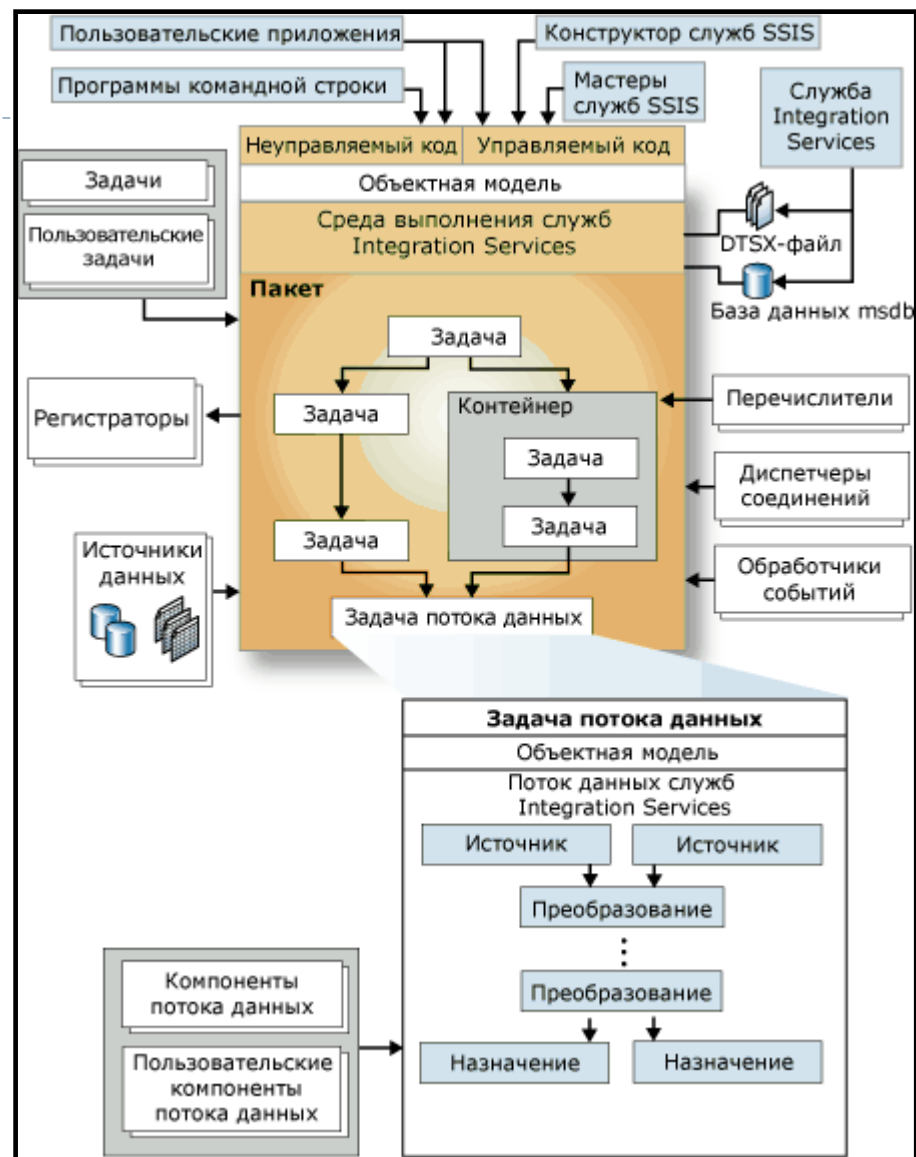
Архитектура SQL Server 2008 Integration Services

Термины

- Источник (и) - Sources
 - Приёмник(и) - Destinations
 - Преобразование данных - (Transformation)
-
- Время исполнения
-
- Пакет (Package)
 - Задача (Task)
 - Буфер (Buffer)
 - Труба (pipeline) потока данных

Конструктор служб SSIS

- Поток управления (Control Flow)
- Поток данных (Data Flow)
- Обработчики событий в пакете и объектов пакета (Event Handlers)
- Просмотр содержимого пакета
- Просмотр выполнения пакета



Типовые сценарии в Integration Services

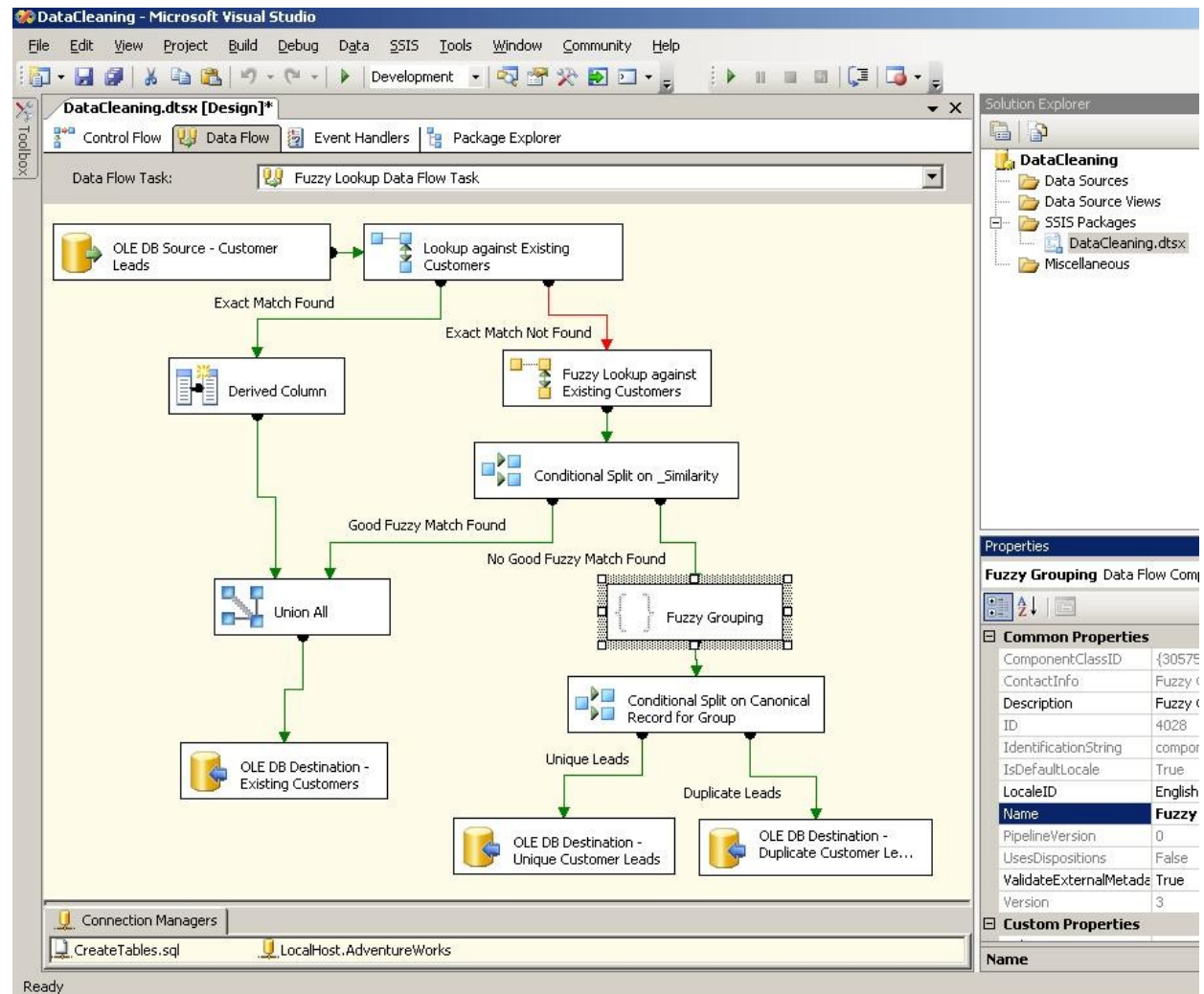
- ▶ Слияние данных из гетерогенных хранилищ данных
- ▶ Очистка, преобразование и стандартизация данных
- ▶ Заполнение хранилищ данных и витрин данных
- ▶ Встраивание бизнес-аналитики в процесс преобразования данных
- ▶ Автоматизация административных функций и загрузки данных

Пример: Очистка данных

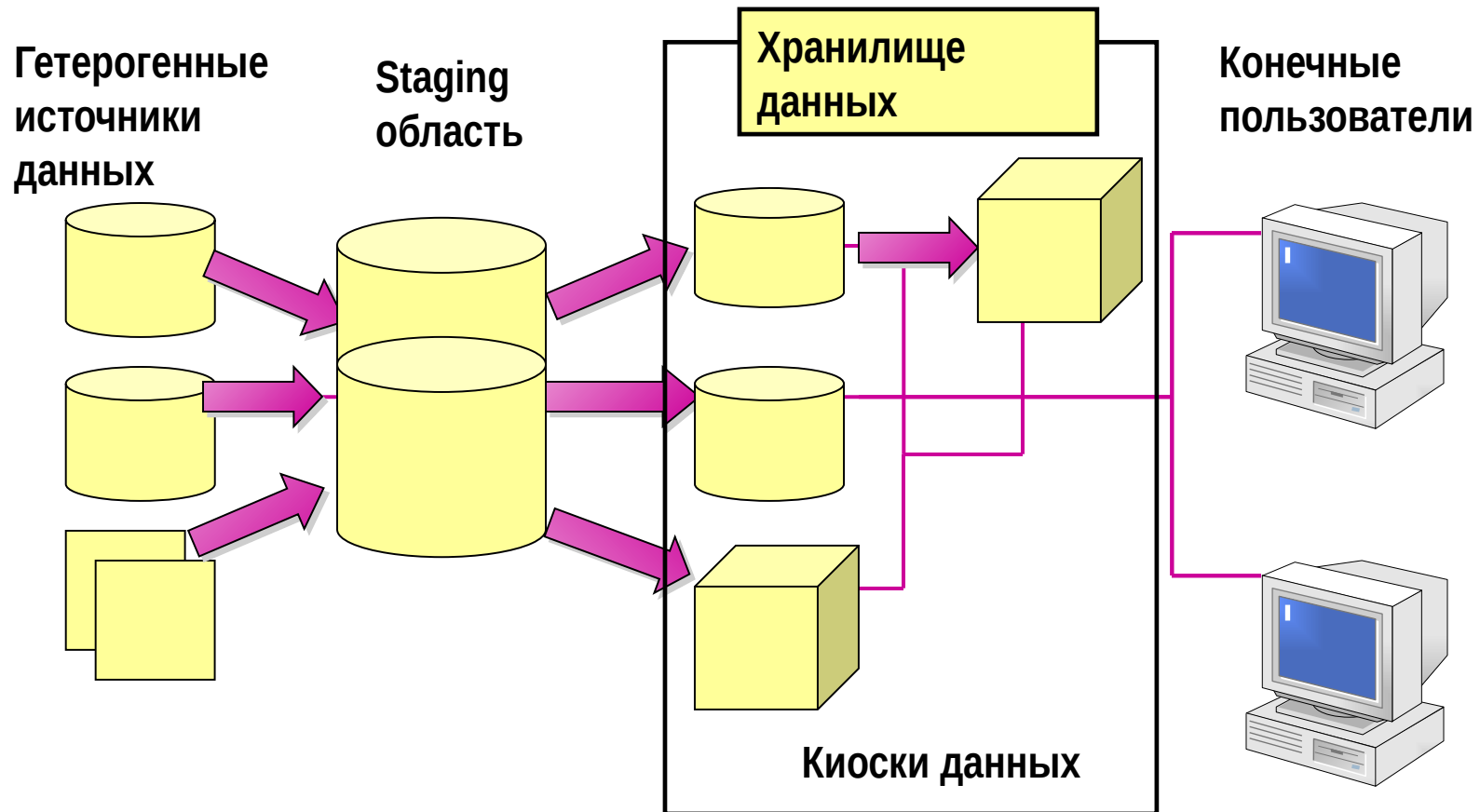
Пакет SSIS
Data Cleaning
Sample из
Integration
Services
Samples.

Fussy Lookup –
нестрогое
соответствие
новых
клиентов
старым
записям

Fussy Grouping –
нечеткий
поиск
фамилий
дубликатов.



Планирование ETL проекта для хранилища данных



Заполнение хранилища данных в SSIS

- Источники и приемники данных
- Оценка и проверка исходных данных
- Промежуточное хранение данных (Staging storage)
- Загрузка в хранилище и витрины данных

Источники и приемники данных

- ▶ Выбрать источники данных (все форматы)
- ▶ Выбрать приемники данных (DW, Data Mart), определить структуру записываемых данных
- ▶ Определить **время** извлечения и записи данных (extraction and load windows),
длительность извлечения и загрузки данных
- ▶ Документировать **диаграмму потока данных**: описать список источников, методов доступа, учетные записи, протоколы, характеристики сети

Промежуточное хранение данных (Staging storage)

В сложных ETL процессах может потребоваться промежуточное хранение данных после чтения перед загрузкой в хранилище:

- Реляционная БД
- Файлы «как есть» - raw (binary) files

После извлечения данных:

- Необходимость быстро освободить источник данных
- Выполнение ETL с заданной контрольной точки без повторного рестарта

Перед загрузкой данных:

- Асинхронное поступление данных, ожидание всех данных
- Фиксируется моментальный снимок данных на заданную дату, возможность получения отчетности по этому снимку данных
- Возможность рестарта с контрольной точки без необходимости выполнять пакет с самого начала
- Возможность провести трансформацию некоторых данных на SQL Server перед окончательной загрузкой в хранилище
- Возможность проверить и удалить невалидные данные или дубликаты после окончания трансформаций перед загрузкой

Загрузка в хранилище и витрины данных

- Загрузка измерений и мер
- Создание первичных и вторичных ключей
- Создание индексов
- Удаление временных таблиц
- Обработка измерений и секций кубов