

ETL : Conceitos e Tecnologias

André Gomes
andregoems@gmail.com
github.com/andregoems
www.linkedin.com/in/andregoems

O que é ETL ?

- **Termo genérico para movimentação e manipulação de dados.**
- **ETL** vem do inglês **Extract Transform Load**
 - **E - Extract** (Lê dado na origem)
 - **T - Transform** (Transforma - Limpa, corrige, agrega , etc.)
 - **L - Load** (Grava no destino **pronto para uso**)
- É um conceito - **NÃO É FERRAMENTA.**
- Processo de ETL popula DW/DM.
- Ferramentas:
 - Informatica PowerCenter .
 - SQL Server Integration Services (SSIS) .
 - Oracle Data Integrator (ODI).
 - IBM InfoSphere DataStage.
 - Talend Data Integration.
 - **Pentaho data integration(PDI).**

ETL - Ferramenta X Programação

ETL com Ferramenta:

- Desenvolvimento mais simples e mais rápido. (Arrastar,soltar e configurar)
- Metadados integrados.
- Agendador interno.
- Linhagem de dados.
- Conectores de várias tecnologias (driver).

ETL com linguagem de programação (Python, Scala, Java, .Net, C# etc..) :

- Flexibilidade ilimitada (Limitada pela linguagem escolhida).
- Teste de unitário

Qual usar ? Depende de:

- Política.
- Cultura.
- Capacidade.
- Projeto.

DEMO



Vendas

id venda	valor venda	id comprador
1	42	4
2	30	2
3	82	2
4	98	2
5	59	2
6	40	2
...
...
...
37	52	4
38	33	2
39	40	2
40	50	2

compradores

id comprador	nome	Estado	idade	estado civil
1	Fulano silva	SP	47	casado
2	Sicrano silva	PE	33	casado
3	Beltrano silva	MG	19	solteiro
4	Sicrano dos santos	SP	29	casado

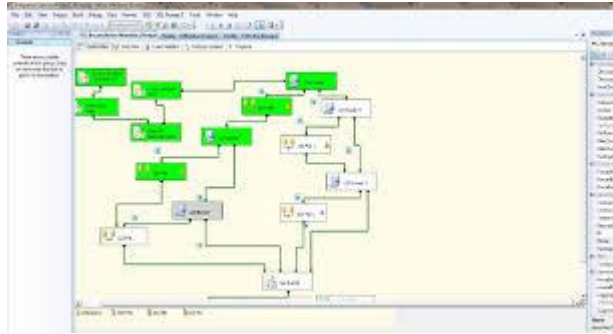
Join (SQL)

Groupby(SQL)

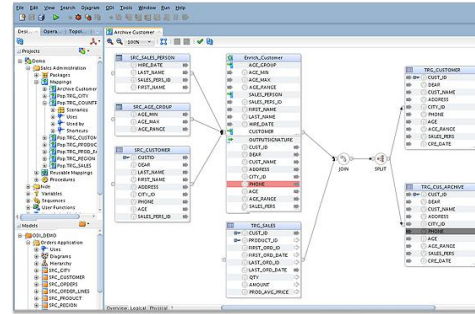
Estado	valor venda
MG	717
PE	1169
SP	491

Ferramenta ETL = Lógica dentro de componente (Caixinha)

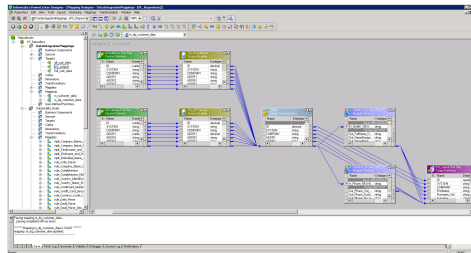
Microsoft SSIS



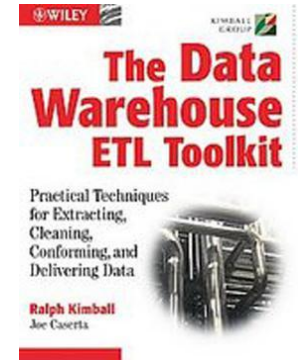
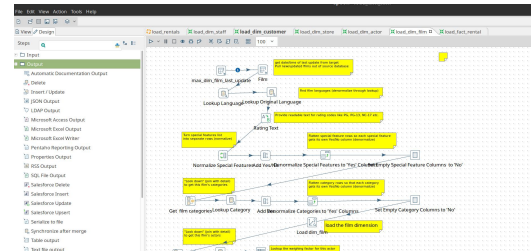
Oracle ODI



Informatica PowerCenter



Pentaho PDI



WILEY

THE KIMBALL GROUP

The Data Warehouse ETL Toolkit

Practical Techniques for Extracting, Cleaning, Conforming, and Delivering Data

Ralph Kimball
Joe Caserta

Sakila Rental Star-Schema

```

    erDiagram
        actor ||--o{ film_actor : "has"
        actor ||--o{ film_text : "has"
        actor ||--o{ staff : "has"
        actor ||--o{ rental : "has"
        actor ||--o{ payment : "has"
        actor ||--o{ category : "has"
        actor ||--o{ film : "has"
        actor ||--o{ store : "has"
        actor ||--o{ address : "has"
        actor ||--o{ city : "has"
        actor ||--o{ country : "has"
        actor ||--o{ customer : "has"

        actor {
            int actor_id PK
            varchar45 first_name
            varchar45 last_name
            timestamp last_update
        }

        film_actor {
            int actor_id FK
            int film_id FK
            timestamp last_update
        }

        film_text {
            int film_id FK
            varchar45 film_title
            timestamp last_update
        }

        staff {
            int staff_id PK
            int inventory_id FK
            varchar45 first_name
            varchar45 last_name
            timestamp last_update
        }

        rental {
            int rental_id PK
            int inventory_id FK
            int customer_id FK
            timestamp last_update
        }

        payment {
            int payment_id PK
            int customer_id FK
            int staff_id FK
            int rental_id FK
            int inventory_id FK
            int amount
            timestamp last_update
        }

        category {
            int category_id PK
            varchar23 name
            timestamp last_update
        }

        film {
            int film_id PK
            varchar45 title
            timestamp last_update
        }

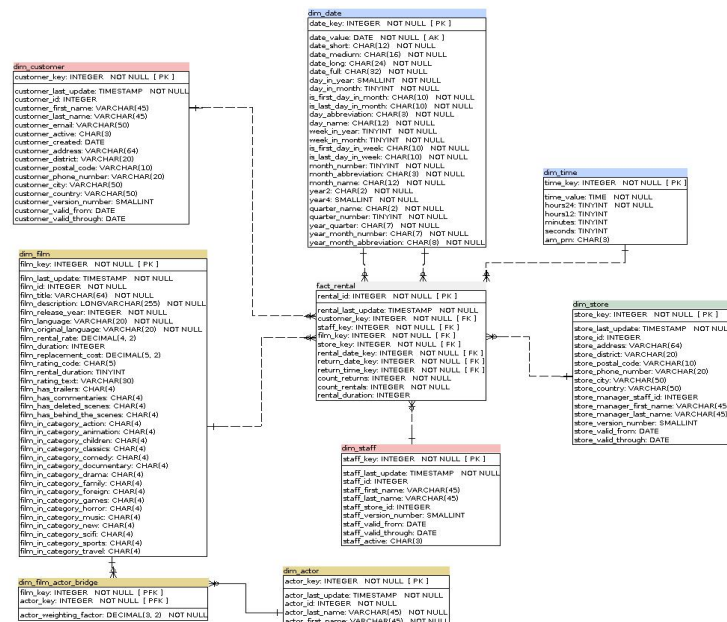
        store {
            int store_id PK
            varchar45 manager_staff_id FK
            timestamp last_update
        }

        address {
            int address_id PK
            varchar45 address
            varchar45 address2
            varchar20 district
            varchar10 postal_code
            varchar10 phone
            timestamp last_update
        }

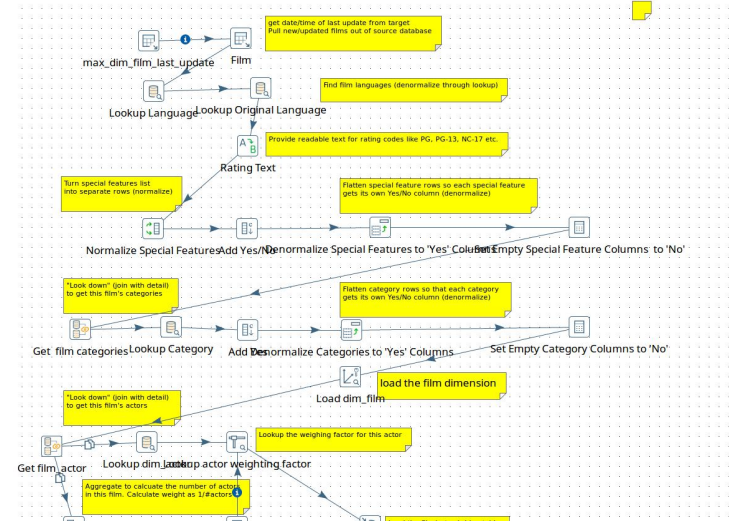
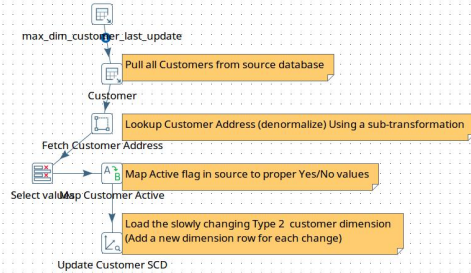
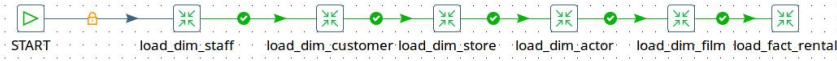
        city {
            int city_id PK
            varchar45 city
            timestamp last_update
        }

        country {
            int country_id PK
            varchar45 country
            timestamp last_update
        }

        customer {
            int customer_id PK
            varchar45 first_name
            varchar45 last_name
            timestamp last_update
        }
  
```

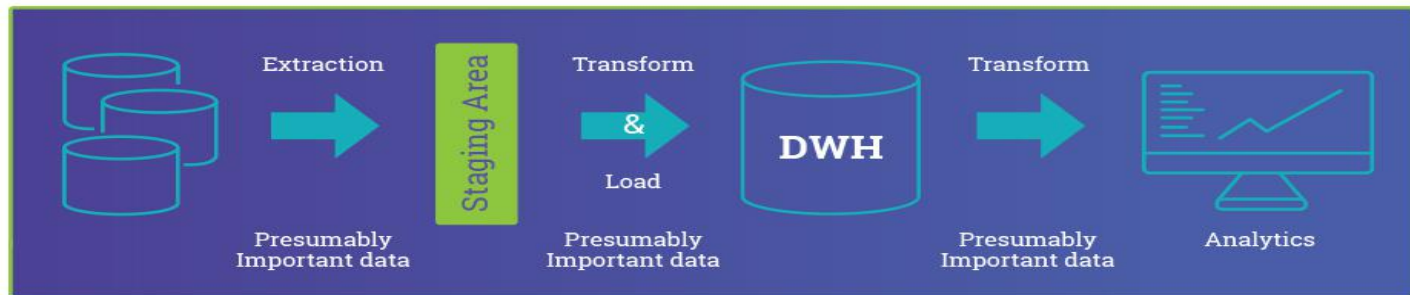


DEMO

ETL x ELT

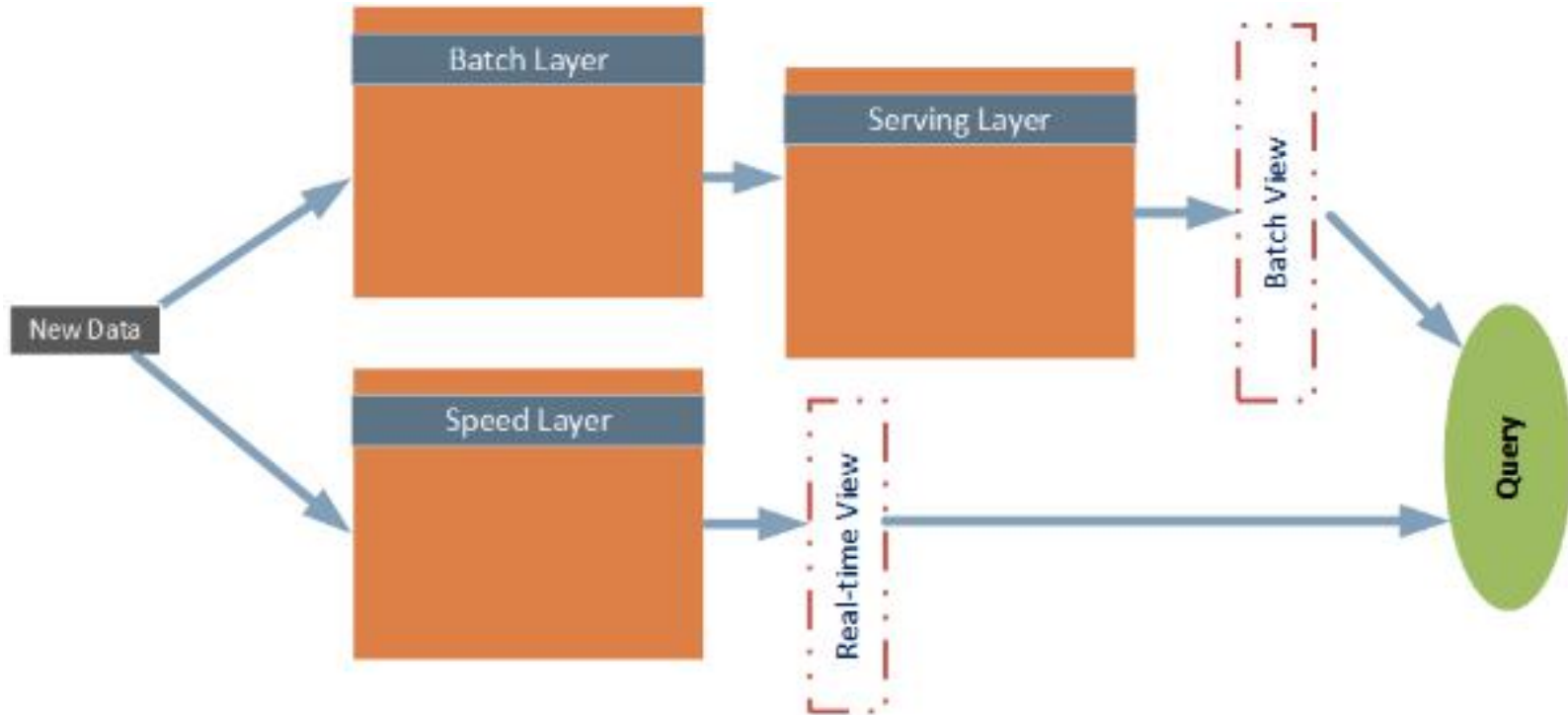
ETL



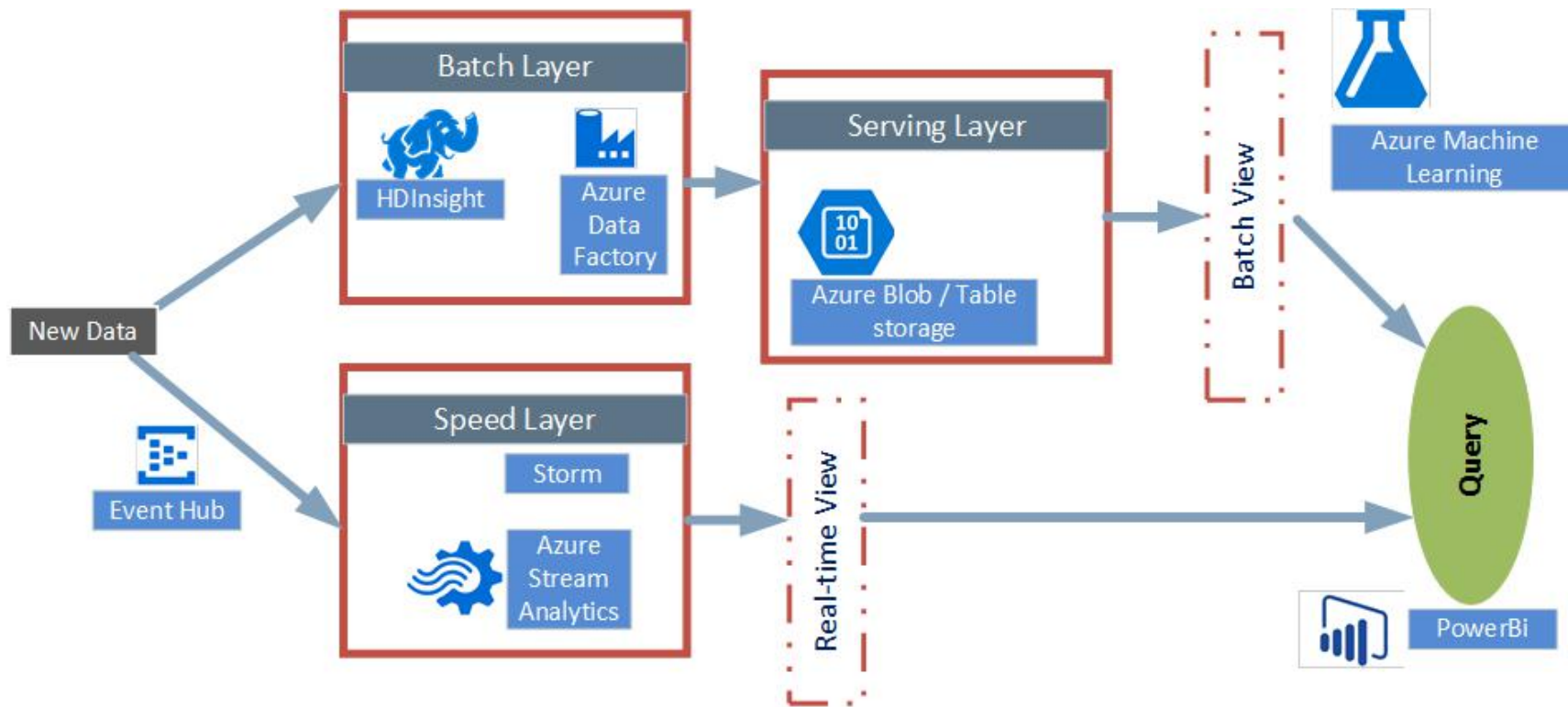
ELT



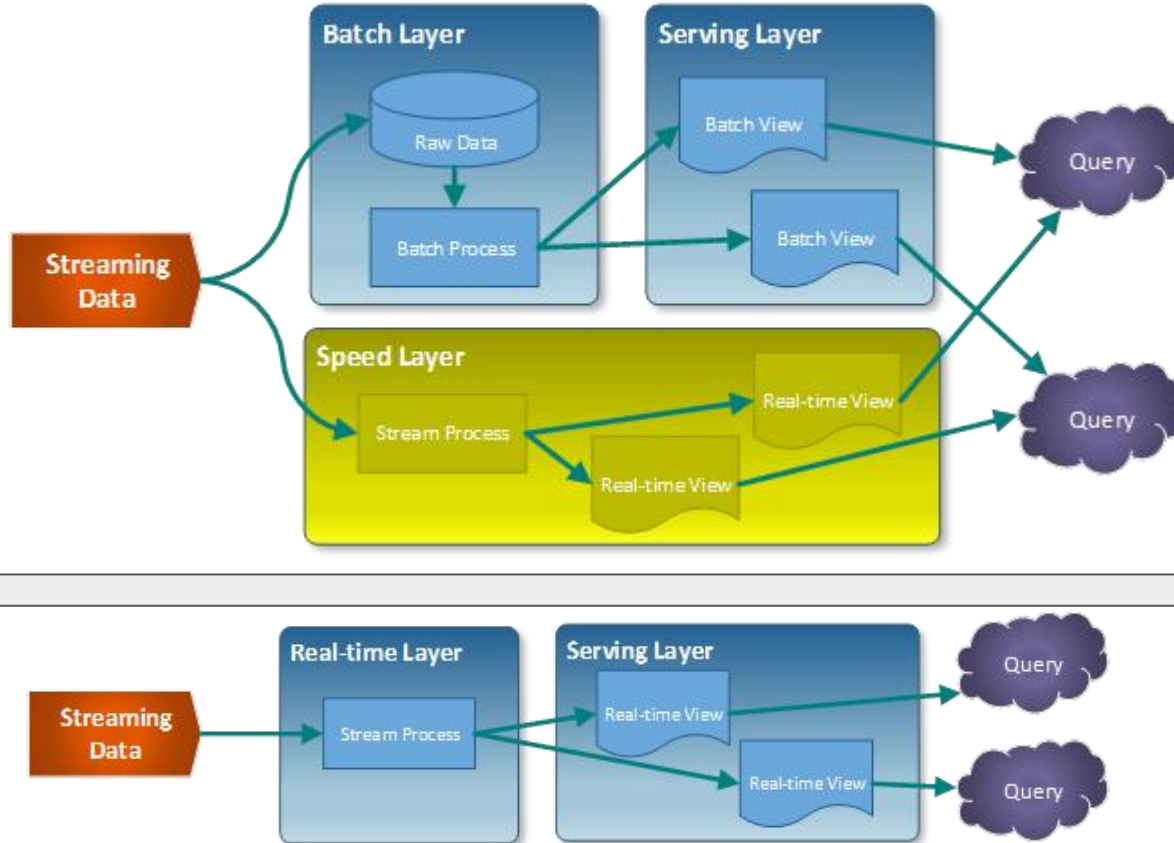
Batch + Streaming = Arquitetura Lambda



Batch + Streaming = Arquitetura Lambda (Azure)



Arquitetura Lambda X Arquitetura Kappa



Big data + Nuvem + Arquitetura = Novas Ferramentas/Frameworks ETL

- Melhor Integração com nuvem ou já estão em nuvem.
- Streaming First.
- Serveless.
- Flexibilidade para novas arquiteturas.
- Volume de dados.
- Integradas com Hadoop, Kafka, Spark, python, R,Hive, etc.

Novas Ferramentas

- Microsoft Data Factory -- <https://azure.microsoft.com/pt-br/services/data-factory/>
- AWS Glue -- <https://aws.amazon.com/glue/>
- Google Cloud Dataflow -- <https://cloud.google.com/dataflow/>
- StreamSets -- <https://streamsets.com/>
- Confluent -- <https://www.confluent.io/product/confluent-platform/>
- Alooma -- <https://www.alooma.com/platform>
- fivetran -- <https://fivetran.com/>
- matillion -- <https://www.matillion.com/platform/>
- snaplogic -- <https://www.snaplogic.com/products>
- Apache Airflow -- <https://airflow.apache.org/>
- Apache NiFi -- <https://nifi.apache.org/>

O que é ETL ?

- Em contexto de BI/DW/DM quase sempre é feito com alguma ferramenta.
- Software para dados = ETL ? (Excel, PowerBI).
- Novos requisitos (velocidade, variedade, volume) = Exige outra maneira de trabalhar com dados.
- Nuvem, Big Data, Arquitetura = Mudança.
- Fabricantes = integração.
- ETL = ELT ? confusão !?!?
- **Depende do seu contexto.**

**Termo genérico para movimentação e
manipulação de dados.**



André Gomes
andregoems@gmail.com
github.com/andregoems
www.linkedin.com/in/andregoems

Obrigado !

André Gomes
andregoems@gmail.com
github.com/andregoems
www.linkedin.com/in/andregoems