



Universidade Federal de Pernambuco
Cin - Centro de Informática

Pós graduação em Ciência da Computação

ETL4NoSQL: Um framework de ETL para BDs NoSQL

Carine Calixto Agüena

Dissertação de Mestrado

Recife
<DATA DA DEFESA>

Universidade Federal de Pernambuco
Cin - Centro de Informática

Carine Calixto Aguená

ETL4NoSQL: Um framework de ETL para BDs NoSQL

*Trabalho apresentado ao Programa de Pós graduação em
Ciência da Computação do Cin - Centro de Informática da
Universidade Federal de Pernambuco como requisito par-
cial para obtenção do grau de Mestre em Ciência da Com-
putação.*

Orientador: Valéria Cesário Times

Recife
<DATA DA DEFESA>

<DIGITE A DEDICATÒRIA AQUI>

Agradecimentos

<DIGITE OS AGRADECIMENTOS AQUI>

<DIGITE AQUI A CITAÇÃO>
—<AUTOR> (<NOTA>)

Resumo

<DIGITE O RESUMO AQUI>

Palavras-chave: <DIGITE AS PALAVRAS-CHAVE AQUI>

Abstract

Keywords: <DIGITE AS PALAVRAS-CHAVE AQUI>

Sumário

1	Introdução	1
1.1	Contextualização	2
1.2	Motivação	3
1.3	Objetivos	4
1.3.1	Objetivo Específico	4
1.4	Organização do Trabalho	4
2	Fundamentação Teórica	5
2.1	ETL	6
2.2	Data Warehouse	6
2.3	Bancos de Dados NoSQL	6
2.3.1	Banco de dados Orientados à Documentos	6
2.3.2	Banco de dados Famílias de Colunas	6
2.3.3	Banco de dados Baseado em Grafos	7
2.3.4	Banco de dados Chave-Valor	7
2.4	Projeto Conceitual, Lógico e Físico	7
2.4.1	Modelo Conceitual NoSQL	8
2.4.2	Modelo Lógico NoSQL	8
2.5	Frameworks	9
2.6	Trabalhos Correlatos	9
2.6.1	Experiência do Usuário	9
2.6.2	Solução Spotfire	9
3	O Framework ETL4NoSQL	11
3.1	Requisitos de software do ETL4NoSQL	12
3.2	Arquitetura do ETL4NoSQL	13
3.3	Componentes do ETL4NoSQL	14
3.3.1	Componente de Importação	16
3.3.2	Componente de Mapeamento	16
3.3.3	Componente de Mecanismos de ETL	17
3.3.4	Componente de Operações	18
3.4	Considerações finais	18

SUMÁRIO

4	Ambiente de Programação	19
4.1	Implementação	20
4.2	Interfaces de Programação	20
4.2.1	Módulo NoSQL	20
4.2.1.1	Esquema de Dados - SchemaData	21
4.2.1.2	Esquema de Dados Família de Coluna - SchemaDataColumn-Family	21
4.2.1.3	Esquema de Dados Documento - SchemaDataDocument	22
4.2.1.4	Esquema de Dados Grafo - SchemaDataGraph	23
4.2.1.5	Esquema de Dados Chave Valor - SchemaDataKeyValue	23
4.2.1.6	Sintaxe DDL - SyntaxDDL	23
4.2.1.7	Sintaxe DDL Cassandra - SyntaxDDLCassandra	23
4.2.1.8	Sintaxe DML - SyntaxDML	24
4.2.1.9	Sintaxe DML Cassandra - SyntaxDMLCassandra	24
4.2.1.10	Inventário - Inventory	24
4.2.1.11	Fábrica de Inventário - InventoryFactory	25
4.2.1.12	Importação - Import	25
4.2.1.13	Amostra - Sample	25
4.2.1.14	Mapeamento	26
4.2.1.15	Amostra Mapeada	26
4.2.2	Módulo ETL	26
4.3	Considerações Finais	26
5	Estudo Experimental de Software	27
5.1	Objetivos do experimento	28
5.1.1	Objetivo da Medição	28
5.1.2	Objetivos do Estudo	28
5.1.3	Questões	28
5.2	Planejamento	29
5.2.1	Definição das Hipóteses	29
5.2.2	Descrição da instrumentação	29
5.2.3	Métricas	29
5.2.4	Seleção do contexto	30
5.2.5	Seleção dos indivíduos	31
5.2.6	Variáveis	31
5.2.7	Análise Qualitativa	31
5.2.8	Validade	31
5.3	Operação	32
5.4	Descrição	32
5.5	Considerações Finais	32

SUMÁRIO

6	Conclusão	33
6.1	Principais Contribuições	34
6.2	Discussão	34
6.3	Resultados	34
6.4	Trabalhos Futuros	34

Lista de Figuras

3.1	Modelo de Processos do ETL4NoSQL	14
3.2	Diagrama de Atividades do ETL4NoSQL	15
3.3	Arquitetura do Framework ETL4NoSQL	15
4.1	Árvore das Classes do Módulo NoSQL	21
4.2	Classe Schema Data	22
4.3	Classe Schema Data Column Family	22
4.4	Classe Schema Data Document	22
4.5	Classe Schema Data Graph	23
4.6	Classe Schema Data Key Value	23
4.7	Classe Syntax DDL	24
4.8	Classe Syntax DML	24
4.9	Classe Syntax DML	25
4.10	Classe Syntax DML	25

Lista de Tabelas

3.1	Requisitos do ETL4NoSQL	13
5.1	Descrição da Instrumentação	30
5.2	Métricas	30

CAPÍTULO 1

Introdução

Este capítulo contextualiza os principais assuntos abordados neste trabalho, apresenta as motivações, os objetivos gerais e específicos da proposta desta pesquisa, bem como sua justificativa.

1.1 Contextualização

Desde a década de 1970, com a criação do modelo relacional por Edgar Frank Codd, a estrutura de armazenamento adotada por muitos desenvolvedores de sistemas da área de tecnologia da informação tem se baseado no conceito de entidade e relação proposto por Codd. A maioria dos sistemas gerenciadores de banco de dados que possui aceitação no mercado fazem uso desse modelo, por exemplo o MySQL, Oracle e Microsoft SQL Server. Porém, os requisitos para o desenvolvimento de ferramentas de software modernas têm mudado significativamente, especialmente com o aumento das aplicações Web [13]. Este segmento de aplicações exige requisitos com alta escalabilidade e vazão, onde sistemas que utilizam um armazenamento com esquema relacional não conseguem atender satisfatoriamente. Em resposta a isso, novas abordagens de armazenamentos de dados utilizando o termo de NoSQL tornaram-se popular.

O termo NoSQL é constantemente interpretado como "*Not Only SQL*", cujo SQL refere-se a linguagem de manipulação de dados dos gerenciadores de armazenamento de dados relacionais (RDBMS - Relational Database Management System) - Structure Query Language [13]. O grande propósito das abordagens NoSQL é oferecer alternativas onde os esquemas relacionais não apresentam um bom desempenho. Esse termo abrange diferentes tipos de sistemas. Em geral, banco de dados NoSQL usam modelo de dados não-relacionais, com poucas definições de esquema, são executados em clusters e aplicados a alguns bancos de dados recentes como o Cassandra, o Mongo, o Neo4J e o Riak [8].

Muitas empresas coletam e armazenam milhares de gigabytes de dados por dia, no qual a análise desses dados torna-se uma vantagem competitiva no mercado. Por isso, há uma grande necessidade de uma nova arquitetura para o gerenciamento de suporte à decisão que possa alcançar melhor escalabilidade e eficiência [12]. Para auxiliar no processo de gerenciamento de suporte à decisão uma das formas mais utilizadas é a criação de um ambiente data warehousing que é responsável por providenciar informações estratégicas e esquematizadas a respeito do negócio [1].

Segundo a definição de [11], data warehouse (DW) é uma coleção de dados para o processo de gerenciamento de suporte à decisão orientado a assunto, integrado, variante no tempo e não volátil. Os dados de diferentes fontes de sistemas são processados em um data warehouse central através da Extração, Transformação e Carga (ETL) de maneira periódica.

Ferramentas de ETL são sistemas de software responsáveis por extrair dados de diversas fontes, transformar e customizar os dados e inseri-los no data warehouse. Comumente, esses processos são executados periodicamente, onde a otimização do seu tempo de execução torna-se importante [14].

O projeto de ETL consome cerca de 70% dos recursos de implantação de um DW, pois desenvolver esse projeto é crítico e custoso, tendo em vista que gerar dados incorretos pode acarretar em más decisões. Porém, por algum tempo pouca importância foi dada ao processo de ETL pelo fato de ser visto somente como uma atividade de suporte aos projetos de DW. Apenas a partir do ano 2000, a comunidade acadêmica passou a dar mais importância ao tema [5].

Tradicionalmente, o DW é implementado em uma base de dados relacional, onde o dado é armazenado nas tabelas fato e tabelas dimensões, na qual forma um esquema em estrela

[11]. Por isso, é comum que as ferramentas de ETL utilizadas no mercado atualmente só dêem suporte aos esquemas relacionais. Para oferecer suporte aos sistemas que necessitem utilizar um esquema não relacional de BDs NoSQL em DW, a proposta desse trabalho é especificar um framework programável, flexível e integrado para modelagem e execução de processos ETL em BDs NoSQL.

1.2 Motivação

A integração de dados e os processos de ETL são procedimentos cruciais para a criação de data warehouses e sistemas BI (business intelligence). Porém, os sistemas para ETL e integração de dados são tradicionalmente desenvolvidos para dados estruturados em modelos relacionais que representam apenas uma pequena parte dos dados mantidos por muitas empresas [3] [Russum 2007, Pedersen 2009]. Dessa forma, existe uma demanda crescente para integrar os dados não estruturados e semi estruturados em um repositório unificado. Devido a complexidade desses dados, novos desafios estão surgindo quando lidamos com dados heterogêneos e distribuídos no ambiente de integração [Salem, 2012].

Além disso, muitas empresas encontram dificuldades para lidar com as ferramentas ETL disponíveis no mercado. Aprender a lidar com essas ferramentas pode ser muito custoso em termos financeiros e de tempo, e por isso, acabam optando desenvolver os seus processos por meio de uma linguagem de programação de propósito geral [Awad et al., 2011; Muñoz et al., 2009].

Portanto, este trabalho propõe um framework programável para desenvolvimento de sistemas de ETL que possibilita a integração de dados não estruturados e semi estruturados armazenados em bases NoSQL. O framework possui um ambiente integrado para a importação e mapeamento dos dados, além da modelagem e customização dos processos de ETL. Os processos de importação e mapeamento do framework integram dados não estruturados e semi estruturados. Esses processos possibilitam a leitura e manipulação de dados de bases NoSQL, e também o armazenamento desses dados em bases deste tipo, oferecendo uma alternativa não relacional para a construção de DWs.

Uma alternativa para organizar e manipular grandes volumes de dados sem utilizar um modelo relacional e ainda processá-los e armazená-los de maneira distribuída é fazer o uso de BDs NoSQL [4]. Com isso, surge a necessidade de se promover meios para o uso desses BDs em DWs.

As pesquisas presentes na literatura sobre extração de dados em BDs NoSQL mostram que não há uma ferramenta que seja integrada para o uso de BDs NoSQL, as ferramentas existentes no mercado apenas oferecem a possibilidade para alguns SGBDs NoSQL, ficando a cargo da equipe de implantação do projeto de DW todo o trabalho de modelagem e programação ao se utilizar BDs NoSQL.

[5] aponta em sua pesquisa que muitas empresas evitam ferramentas de ETL disponíveis no mercado, e adotam o desenvolvimento dos processos a partir de uma linguagem de programação de propósito geral, pelo fato dessas ferramentas terem uma longa curva de aprendizagem e grande complexidade no seu uso.

O aumento do uso de banco de dados com esquemas não relacionais baseados no para-

digma NoSQL e a falta de uma ferramenta programável, flexível e integrada, independente de plataforma que dê suporte à extração, transformação e carga em data warehouses para esses esquemas é a grande motivação deste trabalho.

Dessa forma, encontrar uma solução que seja programável, flexível e integrada para extração, transformação e carga dos dados em BDs NoSQL é a proposta deste trabalho.

1.3 Objetivos

O objetivo principal desta pesquisa é especificar um framework programável, flexível e integrado para modelagem e execução de processos ETL de BDs NoSQL. Os objetivos específicos são detalhados a seguir.

1.3.1 Objetivo Específico

Este trabalho de dissertação tem como objetivo específico realizar um estudo experimental de software a fim de caracterizar as principais funcionalidades das ferramentas de ETL na manipulação de dados estruturados, semi estruturados e não estruturados. O estudo experimental poderá comparar a ferramenta proposta, suas vantagens e desvantagens, em relação às ferramentas de ETL encontradas na literatura.

1.4 Organização do Trabalho

CAPÍTULO 2

Fundamentação Teórica

Neste capítulo são apresentados os conceitos relacionados ao desenvolvimento desta pesquisa.

Os conceitos de ETL e Data Warehouse (DW), bem como o termo NoSQL e os paradigmas de esquemas não relacionais mais utilizados pela comunidade acadêmica, o Famílias de Colunas, Orientados à Documentos, Chave-Valor e Baseado em Grafos.

Também são detalhadas as definições de modelagem conceitual e lógica para esquemas não relacionais.

2.1 ETL

2.2 Data Warehouse

2.3 Bancos de Dados NoSQL

Consistem em bancos de dados não relacionais projetados para gerenciar grandes volumes de dados e que disponibilizam estruturas e interfaces de acesso simples (Lima; Mello, 2015). Cada paradigma NoSQL possui um esquema de modelagem diferente, nos quais são divididas pela literatura em quatro categorias amplamente usadas: Chave-Valor, Orientado a Documentos, Famílias de Colunas e Baseado em Grafos ([Fowler, 2013], [Kaur; Rani, 2013]).

As principais características dos banco de dados NoSQL são:

- Distribuído:
- Escalabilidade Horizontal
- Construído para grande volume de dados
- BASE ao invés de ACID
- Modelo de dados não relacional
- Sem definições de esquema
- Não suporta SQL

[13]

2.3.1 Banco de dados Orientados à Documentos

Banco de dados orientados a documentos são capazes de armazenar documentos como dado. Esses documentos podem ser em qualquer formato como XML (eXtensible Markup Language), YAML (Yet Another Markup Language), JSON (JavaScript Object Notation), entre outros. Os documentos são agrupados na forma de coleções, comparando com banco de dados relacional as coleções são como tabelas e os documentos como os registros. Porém, a diferença entre eles é que cada registro na tabela do banco relacional tem o mesmo número de campos, enquanto que nos documentos na coleção do banco de dados orientado a documentos podem ter campos completamente diferentes (Kaur; Rani, 2013).

Existem mais de 15 banco de dados orientados a documentos disponíveis e os mais utilizados são MongoDB, CouchDB e o RavenDB (Kaur; Rani, 2013).

2.3.2 Banco de dados Famílias de Colunas

Banco de dados baseados em Famílias de Colunas são desenvolvidos para abranger três áreas: número enorme de colunas, a natureza esparsa dos dados e frequentes mudanças no esquema. Os dados em Famílias de colunas são armazenados em colunas de forma contínua, enquanto que em bancos de dados relacionais as linhas é que são contínuas. Essa mudança faz com que operações como agregação, suporte para ad-hoc e consultas dinâmicas se tornem mais eficientes (Kaur; Rani, 2013).

A maioria dos bancos de dados baseados em Famílias de Colunas são também compatíveis

com o framework MapReduce, no qual acelera o processamento de enorme volume de dados pela distribuição do problema em um grande número de sistemas. Os bancos de dados de Família de Colunas open-source mais populares são Hypertable, HBase e Cassandra (Kaur; Rani, 2013).

2.3.3 Banco de dados Baseado em Grafos

Bancos de dados baseado em Grafos são como uma estrutura de rede contendo nós e arestas, onde as arestas interligam os nós representando a relação entre eles. Comparando com o modelo Entidade-Relacionamento, o nó corresponde à entidade, a propriedade do nó à um atributo, a relação entre as entidades ao relacionamento entre os nós. Nos bancos de dados relacionais as consultas requerem atributos de mais de uma tabela resultando numa operação de junção, por outro lado, bancos de dados baseado em Grafos são desenvolvidos para encontrar relações dentro de uma enorme quantidade de dados rapidamente, tendo em vista que não é preciso fazer junções, ao invés disso, ele fornece indexação livre de adjacência (Kaur; Rani, 2013).

2.3.4 Banco de dados Chave-Valor

Em Bancos de dados Chave-Valor os dados são organizados como uma associação de vetores de entrada consistindo em pares de chave-valor. Cada chave é única e é usada para recuperar os valores associados a ele. Esses bancos de dados podem ser visualizados como um banco de dados relacional contendo múltiplas linhas e apenas duas colunas: chave e valor. Buscas baseadas em chaves resultam num baixo tempo de execução, além disso, os valores podem ser qualquer coisa como objetos, hashes, entre outros (Kaur; Rani, 2013).

Os bancos de dados Chave-Valor mais populares são Riak, Voldemort e Redis (Kaur; Rani, 2013).

2.4 Projeto Conceitual, Lógico e Físico

Tradicionalmente um projeto de banco de dados é modelado em três fases denominadas conceitual, lógica e física. O projeto conceitual consiste em apresentar um esquema expressivo que modele os dados de um determinado domínio de informação, enquanto que o projeto lógico transforma um esquema conceitual em algo que se aproxima de um modelo de implementação física do banco de dados. Em projetos de banco de dados NoSQL, há poucos trabalhos que abordam uma metodologia para esquemas lógicos baseados em modelagens conceituais (Lima; Mello, 2015).

Dessa forma, esta seção visa aprofundar o tema a respeito de projeto conceitual e lógico em banco de dados NoSQL.

2.4.1 Modelo Conceitual NoSQL

Em bancos de dados relacionais, o modelo conceitual mais utilizado na literatura é o modelo ER (Entidade-Relacionamento) (Fowler, 2013). Contudo, bancos de dados NoSQL necessitam de um modelo conceitual que atenda às suas características.

O desenvolvimento de banco de dados para sistemas NoSQL é usualmente baseado nas melhores práticas, nas quais são especificamente relacionadas ao sistema desenvolvido, com nenhuma metodologia sistematizada (Bugiotti; Cabibbo; Atzeni, 2014). Por isso, Bugiotti et al (2014) desenvolveu uma abordagem baseada no NoAM (NoSQL Abstract Model). Esta abordagem observa que vários sistemas NoSQL compartilham de características de modelagem similares. Uma importante observação é que sistemas NoSQL oferecem operações de acesso aos dados de forma eficiente, atômica e escalável nas unidades de acesso aos dados em uma certa granularidade. Uma representação errada pode levar a incapacidade de garantir a atomicidade das operações importantes e o desempenho pode piorar dependendo magnitude da aplicação.

A metodologia de Bugiotti et al (2014) procede com a identificação dos agregados, onde cada agregado é um grupo de objetos relacionados que podem ser acessados e/ou manipulados juntos. Essa atividade é importante para suportar escalabilidade e consistência. O modelo conceitual desenvolvido por Bugiotti et al (2014) segue o modelo padrão do DDD (Domain-Driven Design), no qual é uma metodologia muito utilizada em orientação à objeto. Dessa forma, para o modelo conceitual NoSQL, é utilizado o diagrama de classe conceitual da UML, definindo as entidades, valores dos objetos e relacionamentos da aplicação.

2.4.2 Modelo Lógico NoSQL

O modelo de dados lógico dominante nas últimas décadas tem sido o modelo relacional (Fowler, 2013). Porém, para bancos de dados NoSQL a modelagem relacional não atende as características para representação lógica de seus dados. Para Fowler (2013), cada solução NoSQL possui um modelo diferente, os quais ele dividiu em quatro categorias amplamente usadas na literatura: chave-valor, documento, famílias de colunas e grafos.

Nesse mesmo contexto, Lima (2015) propôs a utilização de esquemas lógicos para NoSQL que utilizam o conceito de agregados. Ele justifica a escolha pelo fato de que a representação lógica baseada em agregados apoia os requisitos típicos dos bancos de dados NoSQL, oferecendo suporte à escalabilidade, consistência e desempenho. O conceito de agregados é um termo da área Domain-Driven Design (DDD), sendo uma coleção de objetos relacionados, aninhados, representando como uma única entidade (Lima;Mello, 2015).

Agregado é um padrão de domínio usado para definir a propriedade e fronteira do objeto. Ele é um grupo de objetos associados que são considerados como uma unidade em relação a alterações de dados. O Agregado é demarcado pela fronteira que separa os objetos de dentro para fora. Cada Agregado tem uma raiz. A raiz é uma Entidade, e ela é o único objeto que é acessível de fora do agregado. A raiz pode guardar referências para qualquer dos objetos agregados, e os outros objetos podem guardar referências uns dos outros, mas um objeto de fora só pode guardar referências do objeto raiz. Se houver outras Entidades dentro da fronteira, a identidade dessas entidades é local, fazendo sentido somente dentro do agregado (Domain-

Driven Design Quickly, 2006).

Fowler (2013), define um agregado como um conjunto de objetos relacionados que são tratados como uma unidade, mais precisamente, é uma unidade de manipulação de dados e gerenciamento de consistência. Ele afirma também que trabalhar com banco de dados orientados a agregados traz uma semântica mais clara, enfocando a unidade da interação com o armazenamento de dados. Contudo, o motivo mais importante para a utilização da modelagem orientada a agregados em bancos de dados NoSQL é que ela auxilia a execução em um cluster. Quando se opera em um cluster é necessário minimizar o número de nós a serem pesquisados na coleta de dados. Assim, ao incluir os agregados é possível dar a informação ao banco de dados sobre quais partes serão manipuladas juntas e no mesmo nó.

Dessa forma, a utilização de um modelo lógico NoSQL baseado em agregados se justifica pelo fato de que o conceito desse modelo possibilita o gerenciamento de consistência, a execução em cluster e uma semântica mais clara.

2.5 Frameworks

2.6 Trabalhos Correlatos

Esta seção aborda os trabalhos que são correlatos a esta pesquisa, bem como descreve como estes trabalhos diferem do realizado por esta pesquisa.

2.6.1 Experiência do Usuário

Uma abordagem apresentado por Tableau permite em sua ferramenta que o usuário possa adicionar scripts Hive que consequentemente são executados em um Hadoop cluster. A saída é então importada para a memória no Tableau. O usuário define quais dados do sistema de arquivo devem ser processados, e o quanto complexo os dados deverão ser comprimidos para o formato que a aplicação suporte. Utilizando o cluster de várias máquinas esse processamento permite que o Tableau obtenha resposta rápida para um grande volume de dados até mesmo quando não há nenhum mecanismo de busca disponível. Obviamente, esta abordagem exige que o usuário tenha um conhecimento avançado, e também que haja uma análise a respeito da estrutura na qual os dados estão armazenados antes da importação.

[13]

2.6.2 Solução Spotfire

O Framework ETL4NoSQL

Neste capítulo são apresentados os conceitos do framework ETL4NoSQL, que consiste numa plataforma de software para desenvolvimento de sistemas de ETL. Os dados de entrada de ETL4NoSQL podem ser de bases de dados não relacionais, e também de bases de dados relacionais, mais especificamente, o framework permitirá a integração das diversas bases de dados NoSQL que pertencem aos quatro paradigmas de NoSQL: Orientada a documentos, Família de Colunas, Chave-Valor e Baseada em Grafos, além da tradicional base de dados relacional.

O ETL4NoSQL oferece um ambiente integrado para modelar processos de ETL e implementar funcionalidades utilizando uma linguagem de programação independente de uma GUI (*Graphical User Interface* - Interface Gráfica do Usuário).

Para a especificação do framework proposto foram definidas as estruturas de dados dos ambientes de origem, destino e da área de processamento de dados e suas respectivas linguagens de manipulação, e também, as principais funcionalidades dos sistemas de ETL, chamados mecanismos de ETL. Para realizar os processos de ETL, por meio de seus mecanismos, foi definido um controlador de operações que é capaz de se comunicar com os ambientes e os mecanismos de ETL.

A seguir, são detalhados os requisitos de software, a arquitetura do sistema e a estrutura dos componentes utilizados no desenvolvimento do framework.

3.1 Requisitos de software do ETL4NoSQL

Requisitos de software são descrições de como o sistema deve se comportar, definidos durante as fases iniciais do desenvolvimento do sistema como uma especificação do que deveria ser implementado (SOMMERVILLE, 1997). Os requisitos podem ser divididos em funcionais e não funcionais, onde o primeiro descrevem o que o sistema deve fazer, ou seja, as transformações a serem realizadas nas entradas de um sistema, a fim de que se produzam saídas, já o outro expressa as características que este software vai apresentar. (SOMMERVILLE e SAWYER, 1997).

O ETL4NoSQL é um framework que tem como principal objetivo auxiliar na criação de processos de ETL ao se utilizar diversas estruturas de armazenamento de dados. Um sistema de software pode ter seus dados armazenados em bases relacionais, que seguem o modelo entidade e relacionamento, ou não relacionais, onde esta possui pouca definição de esquema, não segue um modelo específico e são regularmente chamados de NoSQL. As bases NoSQL possuem quatro paradigmas frequentemente utilizados: Chave-Valor, Família de Colunas, Documentos e Grafos (REF).

As bases de dados relacionais utilizam uma linguagem de gerenciamento de dados padrão conhecida por SQL (Structure Query Language), porém as bases de dados NoSQL não possuem uma linguagem em comum, como as relacionais, cada estrutura de armazenamento possui sua própria linguagem de gerenciamento de dados. Por isso, é essencial que haja um mecanismo que integre a leitura e escrita dos diversos SGBDs NoSQL.

Outra importante características são os processos de ETL que possuem quatro etapas básicas: extração, limpeza/transformação e carga (Kimball and Caserta, 2004). O fluxo do processo de ETL inicia-se com a extração dos dados a partir de uma fonte, que podem ser bases de dados relacionais, bases NoSQL ou arquivos textuais. A partir da extração, os dados passam para uma Área de Processamento de Dados (APD), onde é possível executar processos de limpeza e transformação por meio de mecanismos de junção, filtro, união, agregação e outros. Finalmente, os dados podem ser carregados em estrutura de dados como repositórios analíticos, data warehouses, ou até mesmo em arquivos Linguagem de Marcação Flexível (XML).

Dessa forma, o ETL4NoSQL deve possuir um ambiente que importe os dados dos diversos SGBDs NoSQL, de arquivos textuais, além dos SGBDs relacionais, e que possa fazer a leitura e escrita dos dados permitindo a execução dos processos de ETL. No quadro 3.1 apresentamos os principais requisitos elencados do ETL4NoSQL. Definimos como importante as prioridades que são imprescindíveis para o desenvolvimento e funcionamento do framework, e desejável as funcionalidades que aprimoram o uso do framework, porém não interferem no seu principal objetivo.

O modelo de processo do funcionamento da ferramenta ETL4NoSQL, baseado nas notações da UML 2.0, é representado na figura 3.1. Esse modelo descreve o processamento dos dados nas atividades de identificação dos dados, obtenção das informações para a importação e o mapeamento dos dados para os esquemas desejados, e também, a atividade dos processos de ETL para por fim dar carga dos dados em DWs, repositórios analíticos ou em arquivos XML.

Outro modelo importante para o entendimento do fluxo de processos da ferramenta ETL4NoSQL é o diagrama de atividades, que de acordo com a UML 2.0 tem como objetivo mostrar o fluxo

Tabela 3.1 Requisitos do ETL4NoSQL

Requisito	Prioridade
O sistema deve importar os dados de diversas bases relacionais e não relacionais	Importante
O sistema deve permitir a leitura e escrita dos dados importados	Importante
O sistema deve permitir mapear os dados no modelo relacional	Importante
O sistema deve permitir mapear os dados em quaisquer modelo desejado pelo usuário	Importante
O sistema deve possuir os mecanismos ETL mais conhecidos na literatura	Importante
O sistema deve possibilitar a criação de novos mecanismos ETL desejado pelo usuário	Importante
O sistema deve possuir um ambiente que possibilite a execução dos mecanismos de ETL em operações	Importante
O sistema deve permitir o reutilização dos seus mecanismos para vários cenários	Importante
O sistema deve permitir processamento distribuído	Desejável
O sistema deve permitir a importação de dados a partir de uma nuvem	Desejável

de atividades em um único processo. O diagrama mostra como um atividade depende uma da outra. Na figura 3.2 o diagrama mostra a interação dos componentes ao executar um processo de ETL, onde o estágio inicial é a importação dos dados seguido pelo mapeamento, após a obtenção dos dados necessários é possível a execução dos diversos processos de ETL em uma área de processamento para finalmente os dados serem exportados para base de destino.

3.2 Arquitetura do ETL4NoSQL

Sommerville (2007), define o projeto de arquitetura como um processo criativo em que se tenta organizar o sistema de acordo com os requisitos funcionais e não funcionais. Um estilo de arquitetura é um padrão de organização de sistema (Garlan e Shaw, 1993; Sommerville, 2007), como uma organização cliente-servidor ou uma arquitetura em camadas. Porém, a arquitetura não necessariamente utilizará apenas um estilo, a maioria dos sistemas de médio e grande porte utilizam vários estilos. Para Garlan e Shaw, há três questões a serem definidas na escolha do projeto de arquitetura, a primeira é a escolha da estrutura, cliente-servidor ou em camadas, que permita atender melhor aos requisitos. A segunda questão é a respeito da decomposição dos subsistemas em módulos ou em componentes. E por fim, deve-se tomar a decisão de sobre como a execução dos subsistemas é controlada. A descrição da arquitetura pode ser representada graficamente utilizando modelos informais e notações como a UML (Clements, et al., 2002; Sommerville, 2007).

A arquitetura do ETL4NoSQL, representada graficamente na figura 3.3, é baseada no re-

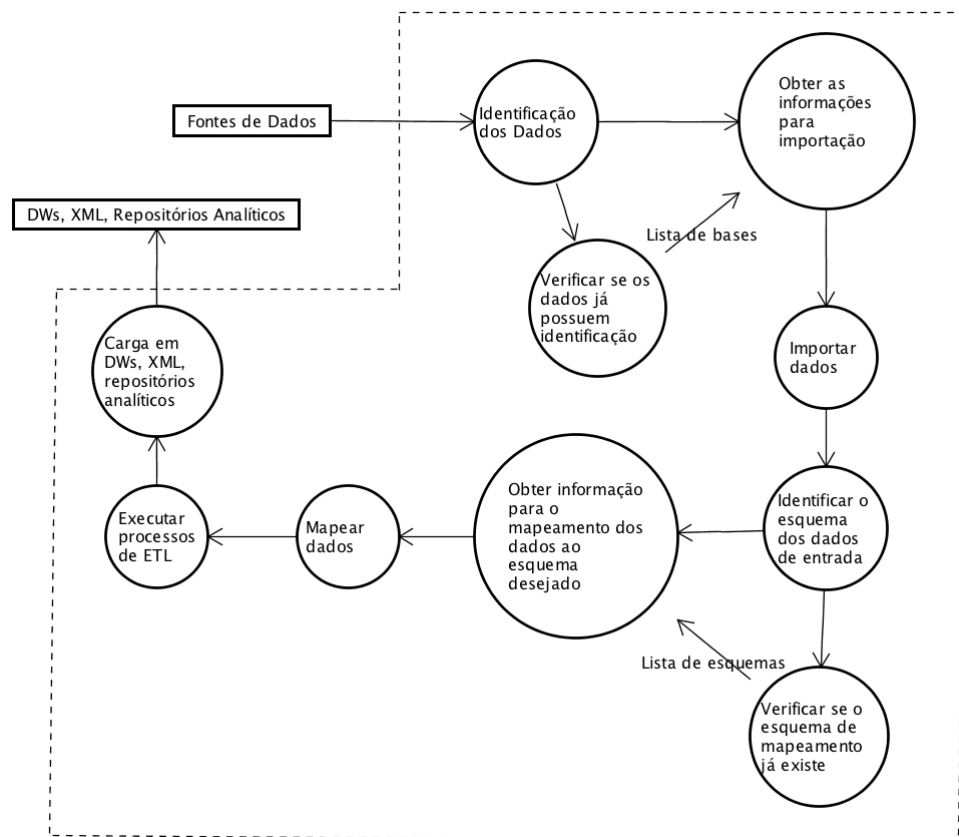


Figura 3.1 Modelo de Processos do ETL4NoSQL

quisito de reutilização. A possibilidade do reuso, reduz o trabalho repetitivo na implementação de componentes e o custo de manutenção (Szyperski, et al.2002), e sua estrutura é em camadas, onde há a camada de sistema e a camada de interface. A camada de sistema lida com todas as operações internas e a camada de interface faz toda a interligação do sistema com o ambiente externo. A decomposição dos subsistemas do ETL4NoSQL é em componentes, pois componentes podem ser subsistemas ou simples objetos que podem ser reusados (Sommerville, 2007). Os componentes que integram o framework e representados na figura 3.3 são os componentes de importação, mapeamento, mecanismos ETL e Operações. Estes componentes serão melhor detalhados na seção seguinte.

3.3 Componentes do ETL4NoSQL

A engenharia de software baseada em componentes é uma abordagem baseada em reuso para desenvolvimento de sistemas de software, ela envolve o processo de definição, implementação e integração ou composição de componentes independentes não firmemente acoplados ao sistema. Os componentes são independentes, ou seja, não interferem na operação uns dos outros e se comunicam por meio de interfaces bem definidas, os detalhes de implementação

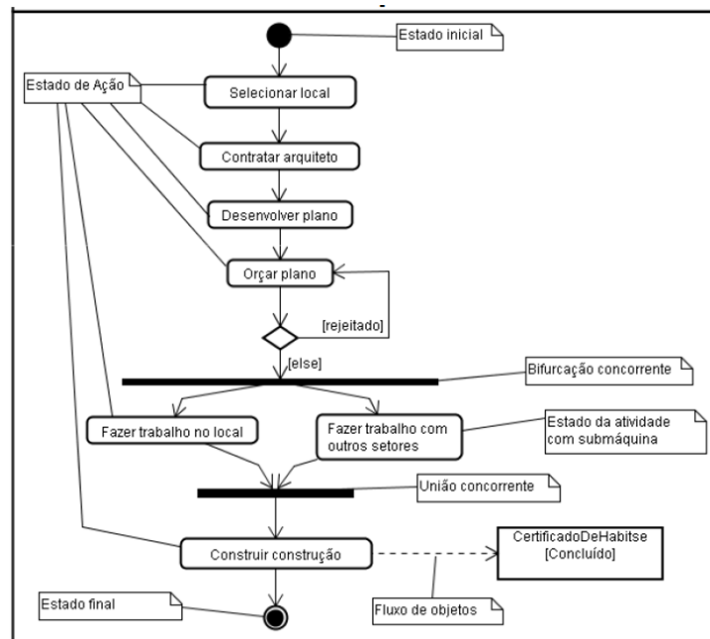


Figura 3.2 Diagrama de Atividades do ETL4NoSQL

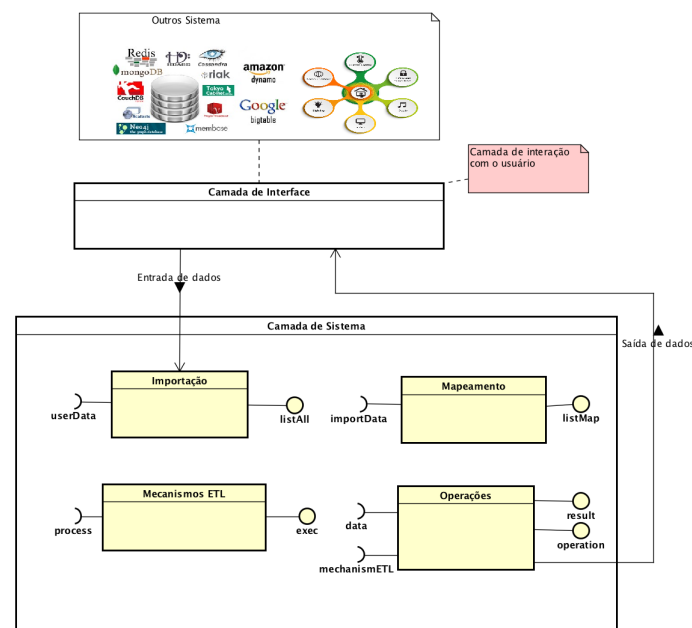


Figura 3.3 Arquitetura do Framework ETL4NoSQL

são ocultados, de forma que as alterações de implementação não afetam o restante do sistema (Sommerville, 2007). Segundo [15], componentes são uma parte do sistema de software que podem ser identificados e reutilizados, onde descrevem ou executam funções específicas e possuem interfaces claras, documentação apropriada e a possibilidade de reuso bem definida.

Ainda de acordo com o autor, um componente deve ser autocontido, identificável, funcional, possuir uma interface, ser documentado e ter uma condição de reuso.

De acordo com os requisitos do ETL4NoSQL, foi possível identificar quatro importantes funcionalidades que podem ser definidas como componentes do sistema, a funcionalidade de importação, mapeamento de dados, mecanismos dos processos de ETL e o controlador de operações. Os componentes do ETL4NoSQL e suas características são apresentados nas seções seguintes, seguindo as características de componentes adotadas por [9].

3.3.1 Componente de Importação

Um dos objetivos do framework ETL4NoSQL é possibilitar a integração de vários tipos de esquemas de dados, relacionais ou não relacionais, presentes nos sistemas modernos. Para isso, a ferramenta deve permitir a leitura e escrita dos diversos SGBDs existentes que aplicam esses esquemas. A solução encontrada para isso foi desenvolver um componente programável que possibilite a importação dos dados por meio de inserção de parâmetros em linha de comando. Este componente, por ser criado utilizando o paradigma de orientação a objetos, permite também sua extensão, por meio de especialização, para que atenda a especificidade de cada cenário. As características do componente são apresentadas a seguir.

- a) Interface: Componente responsável pela importação dos dados da base de origem.
- b) Nomeação: Import.
- c) Metadados: Este componente contém as informações da base de origem como a linguagem de manipulação de dados e meios para estabelecer a conexão com a base, requer uma interação com a interface para o usuário disponibilizar as informações e fornece os dados importados para outros componentes.
- d) Interoperabilidade: Oferece comunicação com outros componentes por meio dos métodos listAll e userData.
- e) Customização: Este componente permite customizar as formas de apresentar os dados importados, de acordo com a necessidade de cada sistema.
- f) Suporte a evolução: Possibilita o suporte aos métodos de acordo com as mudanças de conexões e manipulações de bases de dados futuras.
- g) Empacotamento e utilização: Os métodos são encapsulados e podem ser utilizados pela importação de sua classe e a interface com o usuário é por meio de linha de comando.

3.3.2 Componente de Mapeamento

Para viabilizar a organização dos dados em vários tipos de esquemas desejáveis pelo usuário o ETL4NoSQL oferece o componente de mapeamento. Este componente permite definir o esquema dos dados de acordo com a necessidade da aplicação almejada pelo usuário. Por meio

de parâmetros de inserção em linha de comando é possível utilizar os esquemas de dados pré-definidos pelo componente, mas também, por utilizar o paradigma de orientação a objetos e as características de reusabilidade dos componentes, é possível especializar e customizar os esquemas conforme a conveniência do usuário.

- a) Interface: Componente responsável por gerar o mapeamento dos dados oferecidos pelo componente de importação para um esquema relacional.
- b) Nomeação: Map.
- c) Metadados: Este componente requer os dados de uma base de dados para efetuar o mapeamento.
- d) Interoperabilidade: Oferece comunicação com outros componentes por meio dos métodos `importData` e `listMap`.
- e) Customização: É possível customizar as regras de mapeamento para outros esquemas de dados.
- f) Suporte a evolução: Possibilita o suporte aos métodos de acordo com a necessidade de alterar os esquemas dos dados.
- g) Empacotamento e utilização: Os métodos são encapsulados e podem ser utilizados pela importação de sua classe e a interface com o usuário é por meio de linha de comando.

3.3.3 Componente de Mecanismos de ETL

O ETL4NoSQL é um framework de ETL que possibilita a integração de várias estruturas de dados, por isso ele deve apresentar mecanismos que viabilizem as principais operações de ETL conhecidas pela literatura. Dessa forma, para disponibilizar as operações de ETL, o ETL4NoSQL possui um componente de mecanismos de ETL que permite executar processos de ETL como extração, limpeza/transformação e carga de dados. Além das operações básicas de ETL, o componente permite a especialização e criação de mecanismos permitindo a customização das operações de ETL conforme a necessidade do usuário.

- a) Interface: Componente que contém métodos que realizam as principais operações de ETL presentes na literatura.
- b) Nomeação: MechanismETL.
- c) Metadados: Este componente requer dados de controle para realizar as operações por meio de seus métodos.
- d) Interoperabilidade: Oferece comunicação com outros componentes por meio dos métodos `exec` e `process`.
- e) Customização: É possível customizar e criar mecanismos de acordo com a necessidade de cada processo de ETL.

- f) Suporte a evolução: Deve possibilitar o suporte aos métodos de acordo com a necessidade de alterar os esquemas dos dados.
- g) Empacotamento e utilização: Os métodos deverão ser encapsulados e poderão ser utilizados pela importação de sua classe e a interface com o usuário será por meio de linha de comando.

3.3.4 Componente de Operações

Para proporcionar o controle dos processos de ETL executados pelo framework, o ETL4NoSQL possui o componente de operações. Este componente é responsável pelo controle das operações dos processos de ETL, ele assegura a execução dos mecanismos de ETL de acordo com a necessidade do usuário. É possível também, customizar e especializar as operações deste componente.

- a) Interface: Componente responsável por criar e executar processos de ETL.
- b) Nomeação: Componente de Operação.
- c) Metadados: Este componente deverá possibilitar a comunicação com o componente de mecanismos de ETL e deverá criar e executar processos de ETL.
- d) Interoperabilidade: Deve possibilitar a comunicação entre outros componentes.
- e) Customização: É possível customizar os processos de ETL criados.
- f) Suporte a evolução: Deve possibilitar o suporte aos métodos de acordo com a necessidade de alterar os processos.
- g) Empacotamento e utilização: Os métodos deverão ser encapsulados e poderão ser utilizados pela importação de sua classe e a interface com o usuário será por meio de linha de comando.

3.4 Considerações finais

CAPÍTULO 4

Ambiente de Programação

Este capítulo consiste na apresentação do ambiente de programação do ETL4NoSQL. Ele foi desenvolvido utilizando a linguagem de programação orientada a objetos Python. É demonstrado também os aspectos de implementação, as classes de software e as instâncias dos objetos das classes. As classes são interfaces de programação orientada a objetos fundamentais para o modelo de abstração de frameworks, e as classes deste trabalho serão utilizadas para a importação, mapeamento de BDs NoSQL e criação de processos de ETL a partir desses BDs.

4.1 Implementação

A implementação do ETL4NoSQL foi feita utilizando a linguagem de programação orientada a objetos Python. A escolha dessa linguagem justifica-se pelo fato dela utilizar o paradigma de orientação a objetos que é adequada para a implementação dos padrões de projeto desenvolvidos na proposta deste trabalho. Além disso, Python tem uma sintaxe de fácil aprendizado e pode ser usada em diversas áreas, como Web e computação gráfica. Ela é uma linguagem de alto nível interpretada, completamente orientada a objetos e também é um software livre.

Assim, a implementação do framework foi baseada nos princípios do design orientado a objetos de inversão de controle, onde determina que os módulos de alto nível não devem ser dependentes de módulos de baixo nível, e sim, de abstrações, ou seja, os detalhes devem depender das abstrações. Esse princípio sugere que dois módulos não devem ser ligados diretamente, pois devem estar desacoplados com uma camada de abstração entre eles. Para suprir esse princípio, o ETL4NoSQL tem uma classe abstrata para o Esquema de Dados, que utiliza dos mesmos comportamentos, para as diversas variações de esquemas que os paradigmas NoSQL possui, porém aplicados de acordo com a especificidade de cada um. Outro princípio importante utilizado é o da segregação de interfaces onde os usuários não devem ser forçados a depender de interfaces que não necessitam, deve-se escrever interfaces enxutas com métodos que sejam específicos da interface.

Portanto, o framework foi dividido em interfaces de importação, mapeamento dos BDs NoSQL, mecanismos e operações de processos de ETL. As ferramentas utilizadas para implementação do ETL4NoSQL foram:

- Notebook com sistema operacional MacOS X; processador de 2,5 GHz Intel Core i5; e memória 12 GB 1333 MHz DDR3;
- Python 2.7: Linguagem de programação orientada a objetos.
Disponível em: <https://www.python.org/download/releases/2.7/>;
- LiClipse 3.4.0: plataforma de programação (IDE) open-source.
Disponível em: <http://www.lclipse.com/download.html>;
- SGBD MariaDB versão 10.0.27.
Disponível em: <https://downloads.mariadb.org/mariadb/10.0.27/>;
- SGBD Redis versão 3.2. Disponível em: <https://redis.io/download>;
- SGBD Cassandra 3.0. Disponível em: <http://cassandra.apache.org/download/>

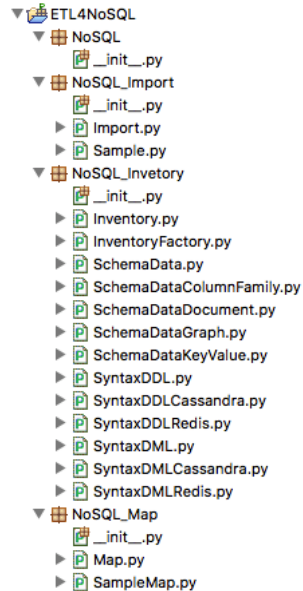
4.2 Interfaces de Programação

4.2.1 Módulo NoSQL

O módulo NoSQL é responsável por lidar com toda a parte que diz respeito ao dados modelados a partir dos paradigmas NoSQL. É neste módulo que serão feitas as importações dos

dados e mapeamentos dos esquemas das bases não relacionais. A árvore de organização das classes do módulo pode ser vista na Figura 4.1.

Figura 4.1 Árvore das Classes do Módulo NoSQL



4.2.1.1 Esquema de Dados - SchemaData

Esquema de dados é uma classe que utiliza o padrão factory method, esse padrão permite que as interfaces criem objetos, porém a responsabilidade de criação fica a cargo da subclasse, as classes que derivam dela são as variações de esquemas existentes dos vários paradigmas de armazenamento de dados presentes na literatura. O trecho do código está ilustrado na Figura 4.2.

O método createSchema é um método abstrato, e é por meio dele que as subclasses implementarão a criação dos seus esquemas de maneira personalizada.

4.2.1.2 Esquema de Dados Família de Coluna - SchemaDataColumnFamily

Classe do tipo ConcreteCreator derivada do Esquema de Dados onde define a criação do esquema das bases sob o paradigma NoSQL Família de coluna. Segundo Nasholm (2012), o esquema de dados do paradigma Família de Coluna é uma coleção de colunas em uma tabela. As linhas são similares às linhas do esquema relacional, exceto que todas as linhas na mesma tabela não necessariamente tem a mesma estrutura. Um valor numa célula, por exemplo, a intersecção de uma linha e uma coluna, é uma sequencia não interpretada de bit. Cada célula é versionada, significando que ela contém múltiplas versões do mesmo dado e que cada versão tem um timestamp atrelada a ela. O trecho do código está ilustrado na Figura 4.3.

Figura 4.2 Classe Schema Data

```

from abc import abstractmethod

class SchemaData:
    def __init__(self, name, describe):
        self.columns = []
        self.name = name
        self.describe = describe
        self.createSchema()

    @abstractmethod
    def createSchema(self):
        pass

    def getSchema(self):
        return "Name Schema: " + self.name + " Describe Schema: " + self.describe

    def putColumn(self, key):
        self.columns.append(key)

    def getColumn(self, key):
        index = self.columns.index(key)
        return self.columns[index]

    def getColumns(self):
        return self.columns

    def popColumn(self, key):
        self.columns.pop(self.columns.index(key))

```

Figura 4.3 Classe Schema Data Column Family

```

from NoSQL_Invetory.SchemaData import SchemaData

class SchemaDataColumnFamily(SchemaData):

    def createSchema(self):
        key = raw_input("Insert key value: ")
        timestamp = raw_input("Insert timestamp ")
        self.putColumn([key, timestamp])

```

4.2.1.3 Esquema de Dados Documento - SchemaDataDocument

Esquema de Dados Documento é uma subclasse do Esquema de Dados do tipo Concrete-Creator, ela define o esquema das bases sob o paradigma NoSQL Orientado a documento. Conforme Nasholm (2012), um documento é geralmente um conjunto de campos onde o campo é um par chave-valor. Chaves são strings atômicas ou sequência de bits, e valores também são atômicos, por exemplo, inteiros ou strings, ou complexos, por exemplo, listas, mapas, entre outros. Um armazenamento de dados de documentos pode armazenar muitos documentos ou até mesmo muitas coleções de documentos. O trecho do código está ilustrado na Figura 4.4.

Figura 4.4 Classe Schema Data Document

```

from NoSQL_Invetory.SchemaData import SchemaData

class SchemaDataDocument(SchemaData):

    def createSchema(self):
        key = raw_input("Insert key value: ")
        document = raw_input("Insert document ")
        self.putColumn([key, document])

```

4.2.1.4 Esquema de Dados Grafo - SchemaDataGraph

Subclasse de SchemaData, define o esquema das bases sob o paradigma NoSQL Baseada em Grafos. De acordo com Nasholm (2012), bases de dados baseada em grafos é estruturada em grafos matemáticos. Um grafo $G=(V, E)$ geralmente consiste em um conjunto de vértices V e um conjunto de arestas E . Uma aresta $e \in E$ é um parte de vértices $(v1, v2) \in V \times V$. Se o grafo é direto esses pares são ordenados. Os vértices do grafo são chamados de nós, e as arestas de relações. Cada nó contém um conjunto de propriedades. O trecho do código está ilustrado na Figura 4.5.

Figura 4.5 Classe Schema Data Graph

```
from NoSQL_Invetory.SchemaData import SchemaData

class SchemaDataGraph(SchemaData):

    def createSchema(self):
        node = raw_input("Insert node value: ")
        edge = raw_input("Insert edge value: ")
        vertice = raw_input("Insert vertice value:")
        self.putColumn([node, edge, vertice])
```

4.2.1.5 Esquema de Dados Chave Valor - SchemaDataKeyValue

Subclasse derivada do SchemaData onde define o esquema das bases sob o paradigma NoSQL Chave-Valor. Para Nasholm (2012), o modelo de dados Chave Valor é baseado na abstração de dados do tipo Map. Ele contém a coleção de pares chave-valor onde todas as chaves são únicas. O trecho do código está ilustrado na Figura 4.6.

Figura 4.6 Classe Schema Data Key Value

```
from NoSQL_Invetory.SchemaData import SchemaData

class SchemaDataKeyValue(SchemaData):

    def createSchema(self):
        key = raw_input("Insert key value: ")
        self.putColumn(key)
```

4.2.1.6 Sintaxe DDL - SyntaxDDL

Classe que utiliza o padrão factory method e define o comportamento da linguagem de definição de dados para cada tipo de SGBD e esquema de dados a serem criados, alterados ou excluídos. O trecho do código está ilustrado na Figura 4.7.

4.2.1.7 Sintaxe DDL Cassandra - SyntaxDDLCassandra

Subclasse de Sintaxe DDL onde define a linguagem de definição de dados do SGBD Cassandra para criação, alteração e exclusão de esquemas com dados oriundos do SGBD Cassan-

Figura 4.7 Classe Syntax DDL

```
from abc import abstractmethod

class SyntaxDDL:

    @abstractmethod
    def createSyntaxDDL(self):
        pass

    @abstractmethod
    def alterSyntaxDDL(self):
        pass

    @abstractmethod
    def dropSyntaxDDL(self):
        pass
```

dra.

4.2.1.8 Sintaxe DML - SyntaxDML

Classe que utiliza o padrão factory method e define o comportamento da linguagem de manipulação de dados para cada tipo de SGBD e esquema de dados a serem manipulados. O trecho do código está ilustrado na Figura 4.8.

Figura 4.8 Classe Syntax DML

```
from abc import abstractmethod

class SyntaxDML:

    @abstractmethod
    def selectSyntaxDML(self):
        pass

    @abstractmethod
    def updateSyntaxDML(self):
        pass

    @abstractmethod
    def insertSyntaxDML(self):
        pass

    @abstractmethod
    def deleteSyntaxDML(self):
        pass
```

4.2.1.9 Sintaxe DML Cassandra - SyntaxDMLCassandra

Subclasse de Sintaxe DML onde define a linguagem de manipulação de dados do SGBD Cassandra para a manipulação dos dados oriundos do SGBD Cassandra. Por meio dessa classe é possível buscar os dados da base de origem Cassandra.

4.2.1.10 Inventário - Inventory

Classe que define um inventário. Um inventário possui nome, descrição, dados de conexão e um conjunto de esquemas. O trecho do código está ilustrado na Figura 4.9.

Figura 4.9 Classe Syntax DML

```
class Inventory:
    def __init__(self, nameInventory, describeInventory, connection):
        self.schemas = []
        self.nameInventory = nameInventory
        self.describeInventory = describeInventory
        self.connection = connection

    def getInventoryName(self):
        return self.nameInventory

    def getInventoryDescribe(self):
        return self.describeInventory

    def getSchemas(self):
        return self.schemas

    def addSchema(self, schema):
        self.schemas.append(schema)
```

4.2.1.11 Fábrica de Inventário - InventoryFactory

Classe que cria inventários, ela é responsável por manter e remover inventários. O trecho do código está ilustrado na Figura 4.10.

Figura 4.10 Classe Syntax DML

```
from NoSQL_Inventoy.Inventory import Inventory
class InventoryFactory:
    def __init__(self):
        self.inventories = []

    def addInventory(self, schemaData, nameInventory, describeInventory):
        self.inventories.append(Inventory(schemaData, nameInventory, describeInventory))

    def removeInventory(self, inventory):
        self.inventories.remove(inventory.getInventoryName())

    def getInventories(self):
        return self.inventories
```

4.2.1.12 Importação - Import

Classe responsável pela importação dos dados das bases NoSQL. Ela utiliza as classes de sintaxes, esquemas, acessa ao inventário e mantém os dados importados.

4.2.1.13 Amostra - Sample

Classe responsável por criar o esquema da amostra que será importada pela classe de importação. Ela acessa o inventário para utilizar o esquema de dados, sintaxes e dados de conexão para criar a amostra.

4.2.1.14 Mapeamento

4.2.1.15 Amostra Mapeada

4.2.2 Módulo ETL

4.3 Considerações Finais

Estudo Experimental de Software

Este capítulo provê o roteiro de experimentação de software para ferramentas de ETL utilizando dados estruturados, semi estruturados e não estruturados. A Engenharia de Software Experimental tem como objetivo aprimorar métodos, técnicas e ferramentas de Engenharia de Software a partir de métodos experimentais (Isaque Elcio de Souza, TESE - Um sistema de Inf para Geren de Projetos Experimentais em ES). As etapas definidas no processo de experimentação em Engenharia de Software proposto por [Amaral (),Isaque Elcio de Souza, TESE] consiste em etapas de definição, planejamento, operação, interpretação dos dados e empacotamento que serão melhor detalhados nas seções a seguir.

5.1 Objetivos do experimento

O objetivo principal da aplicação deste experimento é definir se a ferramenta proposta por esta pesquisa de dissertação é uma ferramenta adequada para auxiliar no desenvolvimento de processos de ETL em dados estruturados, semi estruturados e não estruturados.

5.1.1 Objetivo da Medição

Tendo como base as ferramentas existentes na literatura, caracterizar:

1. Quais as principais funcionalidades que as ferramentas oferecem:
 - (a) essas funcionalidades manipulam dados estruturados, semi estruturados e não estruturados.
 - (b) essas funcionalidades não manipulam dados estruturados, semi estruturados e não estruturados.
2. Quais funcionalidades podem ser consideradas fundamentais para a produtividade na criação de processos de ETL:
 - (a) quais necessitam manipular dados em grande escala.
 - (b) quais não manipulam grande volume de dados.
3. Quais funcionalidades poderiam aprimorar as ferramentas de ETL.

5.1.2 Objetivos do Estudo

Analisar as ferramentas de ETL para dados estruturados, semi estruturados e não estruturados.

Com o propósito de caracterizar.

Com respeito à intersecção das ferramentas de ETL existente.

Do ponto de vista da literatura.

No contexto de comparativo entre as ferramentas mais conhecidas no mercado atual.

5.1.3 Questões

Q1. Existem funcionalidades listadas pelas ferramentas pesquisadas que não estão presentes na ferramenta proposta?

Métrica: A lista de funcionalidades que não estão presentes na ferramenta proposta.

Q2. Existem funcionalidades oferecidas pela ferramenta proposta que não estão presentes nas ferramentas apresentadas pela literatura?

Métrica: A lista de funcionalidades que não estão presentes nas ferramentas da literatura.

Q3. Existem funcionalidades que não estão presentes na ferramenta proposta e nas ferramentas da literatura que poderiam ser implementadas?

Métrica: A lista de funcionalidades que não estão presentes em nenhuma das ferramentas.

5.2 Planejamento

Na etapa de planejamento são definidas as hipóteses do estudo, a descrição da instrumentação, as métricas, seleção do contexto e dos indivíduos, as variáveis, a análise qualitativa e a validade do experimento. Todas elas serão descritas nas seções seguintes.

5.2.1 Definição das Hipóteses

Hipótese nula (H0): As funcionalidades oferecidas pela ferramenta proposta são similares às funcionalidades oferecidas pelas ferramentas presentes na literatura.

Fp - Funcionalidades da ferramenta proposta

Fl - Funcionalidades das ferramentas da literatura

$$H0: F_l - (F_p \cap F_l) = \emptyset$$

Hipótese alternativa (H1): A lista de funcionalidades oferecidas pela ferramenta proposta é diferente da lista de funcionalidades oferecidas pelas ferramentas presentes na literatura.

Fp - Funcionalidades da ferramenta proposta

Fl - Funcionalidades das ferramentas da literatura

$$H1: F_l - (F_p \cap F_l) \neq \emptyset$$

Hipótese alternativa (H2): A lista de funcionalidades que poderiam ser implementadas é diferente da lista de funcionalidades oferecidas pelas ferramentas na literatura e pela ferramenta proposta.

Fp - Funcionalidades da ferramenta proposta

Fl - Funcionalidades das ferramentas da literatura

Fi - Funcionalidades que poderiam ser implementadas

$$H2: F_i - (F_p \cap F_l \cap F_i) \neq \emptyset$$

5.2.2 Descrição da instrumentação

Para cada funcionalidade presente nas ferramentas apresentada na literatura que são consideradas fundamentais para o funcionamento dos processos de ETL pode ser encontrada no quadro 5.1:

Para cada funcionalidade aplicar teste estatístico Chi-2 para definir:

se pode considerar que essa funcionalidade é fornecida;

se pode considerar que essa funcionalidade é útil;

se pode considerar que essa funcionalidade necessita de melhoria.

Resultado: N funcionalidades com valores (P; M; U) onde P - presença 0 - não presente; 1 - presente; U - utilidade 0 - não é útil; 1 - é útil; melhoria 0 - não necessita melhorar; 1 - necessita melhorar.

5.2.3 Métricas

Na tabela 5.2 são apresentadas as métricas utilizadas neste experimento.

Tabela 5.1 Descrição da Instrumentação

Presença da Funcionalidade (P)	Melhoria da Funcionalidade (M)	Utilidade da Funcionalidade (U)
1. Não está presente	1. Necessita melhorar	1. É útil
2. Está presente parcialmente	2. Não há necessidade de melhoria	2. Não é útil
3. Está presente	3. Pode melhorar, mas não necessidade	3. É parcialmente útil

Tabela 5.2 Métricas

Nº	P	M	U	Descrição da Funcionalidade	Questões
1	0	0	0	Não está presente, não necessita melhorar, não é útil	N/A
2	0	0	1	Não está presente, não necessita melhorar, é útil	Q3
3	0	1	0	Não está presente, necessita melhorar, não é útil	N/A
4	0	1	1	Não está presente, necessita melhorar, é útil	Q3
5	1	0	0	Está presente, não necessita melhorar, não é útil	Q1, Q2
6	1	0	1	Está presente, não necessita melhorar, é útil	Q1, Q2
7	1	1	0	Está presente, necessita melhorar, não é útil	Q1, Q2
8	1	1	1	Está presente, necessita melhorar, é útil	Q1, Q2

5.2.4 Seleção do contexto

De acordo com Travassos (2002), o contexto pode ser caracterizado conforme quatro dimensões:

- o processo: on-line / off-line;
- os participantes: ferramentas de ETL;
- realidade: o problema real / modelado;
- generalidade: específico / geral.

Nosso estudo supõe o processo off-line porque as ferramentas não estão sendo testadas durante todo o tempo da utilização, mas em certo instante. Os participantes são as ferramentas de ETL encontradas na literatura. O estudo é modelado porque as funcionalidades das ferramentas não são caracterizadas durante a resolução do problema real, mas utilizando parâmetros subjetivos (ex. presença, utilidade e necessidade). As funcionalidades da ferramenta proposta são comparadas com as ferramentas presentes na literatura, então, o contexto possui o caráter específico.

5.2.5 Seleção dos indivíduos

Como participantes para o estudo propõe-se utilizar as ferramentas encontradas na literatura. Assume-se que esses indivíduos estão presente em diversos estudos realizados e avaliados no meio acadêmico.

Para a escolha das ferramentas utilizadas neste estudo foi levado em consideração a semelhança da finalidade do uso com a ferramenta proposta. Seria conveniente utilizar para o estudo ferramentas que tem o objetivo de auxiliar processos de ETL em diversas estruturas de dados. Dessa forma, a seleção baseou-se nas características das ferramentas.

5.2.6 Variáveis

Variável independente: A lista de funcionalidades das ferramentas encontradas na literatura.
Variáveis dependentes:

1. A similaridade entre as funcionalidades oferecidas pela ferramenta proposta e as funcionalidades encontradas nas ferramentas da literatura.

Pode receber os valores: Igual, quando todas as funcionalidades tem o valor PMU = { 1, X, X } (métricas 5-8); Diferente, quando todas as funcionalidades tem o valor PMU = { 0, X, X } (métricas 1-4) Similar, quando não se cumprem as condições de "Igual" e "Diferente". O grau de similaridade pode ser avaliado como: $\{ 1, X, X \} / \{ 0, X, X \} + \{ 1, X, X \} * 100\%$

2. A utilidade das funcionalidades similares. Mostra a parte útil das funcionalidades oferecidas pela ferramenta proposta: Parte útil: $\{ 1, X, 1 \} / \{ 1, X, X \} * 100\%$ Parte inútil: $\{ 1, X, 0 \} / \{ 1, X, X \} * 100\%$

3. A melhoria das funcionalidades similares. Mostra a necessidade de melhoria nas funcionalidades oferecidas pela ferramenta proposta: Não necessita melhorar: $\{ 1, 0, X \} / \{ 1, X, X \} * 100\%$ Necessita melhorar: $\{ 1, 0, X \} / \{ 1, X, X \} * 100\%$

5.2.7 Análise Qualitativa

Para analisar a informação referente às funcionalidades não oferecidas na ferramenta proposta, mas que poderiam ser implementadas, propõe-se aplicar a análise qualitativa. Essa análise deve apresentar a lista de funcionalidades presentes nas ferramentas da literatura, que não estão presentes na ferramenta proposta, mas que são consideradas necessárias para facilitar a manipulação de dados estruturados, semi estruturados e não estruturados. Assim, essa análise deve considerar funcionalidades com valor PMU = 0, X, X (métricas 1-4) e a opção "É útil" para "utilidade da funcionalidade".

5.2.8 Validade

Validade interna: como mencionado na parte "Seleção dos indivíduos" para o estudo se propõe a utilizar ferramentas presentes na literatura, que são validadas pelo meio acadêmico. Assim, assume-se que elas são representativas para a população de ferramentas de ETL.

Além disso, para redução da influência dos fatores que não são interesse do nosso estudo e, portanto, para aumento da validade interna do estudo supõe-se utilizar dados das ferramentas mais populares da literatura, cuja a validação já tenha passado por diversas avaliações.

Validade de conclusão: para receber os valores da presença, utilidade e melhorias o teste binomial será utilizado. A verificação de hipótese será feita por meio de simples demonstração de presença ou não de funcionalidades nas listas que representam as variáveis independentes.

Validade de construção: esse estudo está caracterizado pela conformidade das funcionalidades listadas na ferramenta proposta com as funcionalidades reais necessárias para a utilização de ferramentas de ETL. As características das ferramentas de ETL presentes na literatura representa a lista de funcionalidades que uma ferramenta de ETL deve apresentar para mostrar o desempenho adequado do ponto de vista da literatura. As funcionalidades, que tem o maior relacionamento com as ferramentas de ETL do ponto de vista dos pesquisadores, foram escolhidas do conjunto total de funcionalidades das ferramentas de ETL presentes na literatura.

Validade externa: como foi mencionado nas partes "Seleção dos indivíduos" e "Validade interna" os participantes do estudo em geral podem ser considerados representativos para a população da literatura apresentada pela academia. Para avaliação do nível de importância das funcionalidades analisadas foi levada em consideração a frequência que a funcionalidade apareceu nas ferramentas da literatura.

Os materiais utilizados no estudo podem ser considerados representativos e "em tempo" para o problema sob análise, porque se compõem das funcionalidades de ferramentas de ETL presentes na literatura atual.

5.3 Operação

A etapa de operação ocorre após a etapa de planejamento do estudo experimental. Nela é exercido o monitoramento do experimento para garantir que ele esteja ocorrendo conforme foi planejado (Souza Isaque, 2015).

5.4 Descrição

5.5 Considerações Finais

CAPÍTULO 6

Conclusão

6.1 Principais Contribuições

6.2 Discussão

6.3 Resultados

6.4 Trabalhos Futuros

Referências Bibliográficas

- [1] S. Chaudhuri and U. Dayal. An overview of data warehousing and olap technology. *SIG- MODRecord*, v.26, n.1, p.65-74, 1997.
- [2] Max Chevalier, SMohammed El Malki, Arlind Kopliku, Olivier Teste, and Ronan Tournier. Implementing multidimensional data warehouses into nosql. *Proceedings of the 17th International Conference on Enterprise Information Systems*, p.172-183, April 27-30, 2015.
- [3] Jerome Darmont, Omar Boussaid, AJean christian Ralaivao, and Kamel Aouiche. An architecture framework for complex data warehouses. *7th International Conference on Enterprise Information Systems (ICEIS'05), Miami, USA*, pages 370-373, 2005.
- [4] Lucas de Carvalho Scabora. Avaliação do star schema benchmark aplicado a bancos de dados nosql distribuídos e orientados a colunas. Master's thesis, USP - São Carlos, April 2016.
- [5] Mário Sergio da Silva. Um framework para desenvolvimento de sistemas etl. Master's thesis, Universidade Federal de Pernambuco, September 2012.
- [6] FAYAD M. E. and SCHMIDT D. C. Object-oriented application frameworks. *Communications of the ACM* 40 (10): 32-38, 1997.
- [7] FAYAD M. E., SCHMIDT D. C., and JOHNSON R. E. Building application frameworks: object-oriented foundations of framework design. *John Wiley and Sons, New York, USA*, pp. 3-29, 1999.
- [8] Martin Fowler and Pramod J. Sadalage. *NoSQL Essencial Um guia conciso para o mundo emergente de Persistência Poliglota*. Novatec, 2013. First Edition.
- [9] George T. Heineman. *Component-Based Software Engineering: putting pieces together*. Addison-Wesley, 2001.
- [10] Ralph Kimball and J. Caserta. *The Data Warehouse ETL ToolKit*. Robert Ipsen, 2004.
- [11] Ralph Kimball and Margy Ross. *The Data Warehouse Toolkit*. Robert Ipsen, 2002. Second Edition.

- [12] Xiufeng Liu, Christian Thomsen, and Torben Bach Pedersen. Cloudeitl: Scalable dimensional etl for hive. *DB Tech Reports*, July 2013.
- [13] Petter Nasholm. *Extracting Data From NoSQL Databases*. PhD thesis, Chalmers University of Technology, SE-412 96 Goteborg Sweeden, January 2012.
- [14] VASSILIADIS P., SIMITSIS A., and GEORGANTAS P. TERROVITISB M. SKIADOPOU-LOS S. A generic and customizable frameworkfor the design of etl scenarios. *Information Systems - Special issue: The 15th international conference on advanced information systems engineering* 30 (7): 492â525, 2005.
- [15] Johannes Sametinger. *Software Engineering with Reusable Componets*. Springer, 1997.
- [16] BRAGA R. T. V. Um processo para construÃ§Ã£o e instanciaÃ§Ã£o de frameworks baseados em uma linguagem de padrÃµesparaumdomnioespecifico. *Master'sthesis, ICMC/USP, SÃ£oPaulo*, 2003.