

**CENTRO UNIVERSITÁRIO FEEVALE**

**LEANDRO FUSSIGER**

***NLPSEARCH***  
***FRAMEWORK PARA INTEGRAÇÃO DE SISTEMAS DE PLN***  
***E API'S DE MECANISMOS DE BUSCA***

**NOVO HAMBURGO, julho 2007.**

**LEANDRO FUSSIGER**

***NLPSEARCH***  
**FRAMEWORK PARA INTEGRAÇÃO DE SISTEMAS DE PLN**  
**E API'S DE MECANISMOS DE BUSCA**

**Centro Universitário Feevale**  
**Instituto de Ciências Exatas e Tecnológicas**  
**Curso de Ciência da Computação**  
**Trabalho de Conclusão de Curso**

**Professor Orientador: Rodrigo Rafael Villarreal Goulart**

**Novo Hamburgo, julho de 2007.**

## **AGRADECIMENTOS**

Gostaria de agradecer em primeiro lugar a minha família que suportou todo meu mal humor enquanto fazia este trabalho. Em especial ao meu Pai e minha querida Mãe que sempre estiveram ao meu lado, dando força, amor e carinho, ao meu irmão Fabio pela ajuda em alguns pontos deste trabalho.

Aos meus amigos que de alguma forma contribuíram para a realização deste trabalho, em especial ao Diego Cirino Kern e ao Fábio Lagemann, dois grandes amigos que percorreram esta mesma estrada tortuosa, chamada Trabalho de Conclusão de Curso, ou TCC, e que estarão se graduando junto comigo.

Ao meu orientador Rodrigo Rafael Villarreal Goulart, apesar de ser colorado é gente boa, e me auxiliou muito na etapa final deste trabalho.

E por último, porém não com menos importância, a mim, por todo meu empenho e dedicação nesta etapa da minha vida.

## RESUMO

O Processamento da Linguagem Natural (PLN) é a área da computação que focada no estudo da língua e de suas aplicações computacionais. Os pesquisadores de PLN desenvolvem estudos sobre fenômenos diferentes da língua, como por exemplo, a análise sintática, resolução de anáforas e no tratamento de diálogos. Os experimentos envolvem geralmente o uso de uma coleção dos textos chamada corpus. Para criar um corpus é necessário coletar os textos, com base em algum tema, e etiquetá-los de acordo com algum aspecto lingüístico, o que consome recursos humanos e computacionais. Para investigar estas deficiências este trabalho sugere o uso de uma larga coleção de textos, a Internet. Para fazê-lo, *API's* de mecanismos de busca serão utilizadas, como por exemplo, as *API's* do *Google*, *Yahoo* e *Technorati*. Usando conceitos de PLN e ferramentas e busca *API's*, este trabalho propõe a construção de um *framework* de PLN que use a Internet como corpus.

Palavras Chave: PLN. *API's* de Mecanismos de Busca, *Frameworks*.

## **ABSTRACT**

The Natural Language Processing (NLP) is the computation area which focus on the study of the language and its computational applications. NLP researchers develop studies on different language phenomena, as we can see on syntactic analysis, anaphors resolution and the dialogs treatment. The experiments generally involve the use of a collection of texts, called corpus. To create a corpus is necessary to collect texts, based on some subject and label it according to some linguistic aspect, which consumes human and computational resources. To work on this problem, this work suggests the use of a wide collection of texts, the Internet. To make it, Web search APIs (Application Programming Interface) will be used, as the Google, Yahoo and Technorati. Using NLP concepts and tools and search APIs, this work proposes the construction of a NLP framework which uses the web as corpus.

Key words: NLP, Web search APIs, Frameworks

## LISTA DE FIGURAS

<b>Figura 1.1 - Objetos interagem através do envio de mensagens (BEZERRA, 2002)</b>	<b>20</b>
<b>Figura 1.2 - Representação Gráfica de Uma Classe (BEZERRA, 2002)</b>	<b>21</b>
<b>Figura 1.3 - Princípios da orientação a objetos como aplicações do Princípio da Abstração (BEZERRA, 2002)</b>	<b>22</b>
<b>Figura 1.4 - Objeto – Conceito de Encapsulamento</b>	<b>23</b>
<b>Figura 1.5 - Herança para Figuras Gráficas (RUMBAUGH, et al., 1994)</b>	<b>24</b>
<b>Figura 1.6 - Exemplo de Polimorfismo</b>	<b>25</b>
<b>Figura 2.1 – Diagrama de Casos de Uso (GUEDES, 2005)</b>	<b>31</b>
<b>Figura 2.2 – Diagrama de Classes (GUEDES, 2005)</b>	<b>32</b>
<b>Figura 2.3 – Diagrama de Seqüência (GUEDES, 2005)</b>	<b>33</b>
<b>Figura 2.4 – Diagrama de Objetos (GUEDES, 2005)</b>	<b>34</b>
<b>Figura 2.5 – Diagrama de Atividades (GUEDES, 2005)</b>	<b>34</b>
<b>Figura 2.6 – Diagrama de Componentes (GUEDES, 2005)</b>	<b>35</b>
<b>Figura 2.7 – Diagrama de Implantação (GUEDES, 2005)</b>	<b>36</b>
<b>Figura 2.8 – Diagrama de Comunicação (GUEDES, 2005)</b>	<b>37</b>
<b>Figura 2.9 – Gráfico de Máquina de Estados (GUEDES, 2005)</b>	<b>37</b>
<b>Figura 2.10 – Diagrama de Pacotes (GUEDES, 2005)</b>	<b>38</b>
<b>Figura 2.11 – Diagrama de Estrutura Composta (GUEDES, 2005)</b>	<b>39</b>
<b>Figura 2.12 – Diagrama de Visão Geral (GUEDES, 2005)</b>	<b>39</b>
<b>Figura 2.13 – Diagrama Temporal (GUEDES, 2005)</b>	<b>40</b>
<b>Figura 3.1 - MVC</b>	<b>43</b>
<b>Figura 3.2 – Componentes do Desenvolvimento tradicional de aplicações (SILVA, 2000)</b>	<b>48</b>

<b>Figura 3.3 – Componentes do desenvolvimento de aplicações baseadas em frameworks (SILVA, 2000)</b>	<b>49</b>
<b>Figura 6.1 – Tela Inicial do <i>WebJspell</i></b>	<b>73</b>
<b>Figura 6.2 – Análise Morfológica de “O”</b>	<b>74</b>
<b>Figura 6.3 – Análise Morfológica de “Futebol”</b>	<b>74</b>
<b>Figura 6.4 – Análise Morfológica de “deveria”</b>	<b>74</b>
<b>Figura 6.5 – Análise Morfológica de “.”</b>	<b>74</b>
<b>Figura 6.6 – <i>Snowball</i></b>	<b>75</b>
<b>Figura 6.7 – Resultado de Análise do <i>Snowball</i></b>	<b>76</b>
<b>Figura 6.8 – <i>Flip On-Line</i></b>	<b>77</b>
<b>Figura 6.9 – Análise Sintática do <i>Flip On-Line</i></b>	<b>77</b>
<b>Figura 6.10 – Análise Sintática do <i>Flip On-Line</i></b>	<b>78</b>
<b>Figura 6.11 – Tela inicial do <i>JBootCat</i></b>	<b>78</b>
<b>Figura 6.12 – Escolha do local para salvar o novo <i>corpus</i></b>	<b>79</b>
<b>Figura 6.13 – Escolha das palavras</b>	<b>80</b>
<b>Figura 6.14 – Geração das tuplas para pesquisa</b>	<b>80</b>
<b>Figura 6.15 – Escolha do motor de busca</b>	<b>81</b>
<b>Figura 6.16 – Resultado da pesquisa na <i>web</i></b>	<b>81</b>
<b>Figura 6.17 – Processamento da pesquisa pela ferramenta</b>	<b>82</b>
<b>Figura 6.18 – Resumo das ações da ferramenta</b>	<b>82</b>
<b>Figura 7.1 – Diagrama de Classes da <i>Google API</i></b>	<b>85</b>
<b>Figura 7.2 – Diagrama de Seqüência da <i>Google API</i></b>	<b>86</b>
<b>Figura 7.3 – Diagrama de Classes da <i>Yahoo API</i></b>	<b>88</b>
<b>Figura 7.4 – Diagrama de Seqüência da <i>Yahoo API</i></b>	<b>89</b>
<b>Figura 7.5 – Diagrama de Classes da <i>Technorati API</i></b>	<b>90</b>
<b>Figura 7.6 – Diagrama de Seqüência da <i>Technorati API</i></b>	<b>91</b>
<b>Figura 8.1 – Diagrama de Classes do Processo de <i>Stemming</i></b>	<b>105</b>
<b>Figura 8.2 – Diagrama de Seqüência do processo de <i>Stemming</i></b>	<b>106</b>
<b>Figura 8.3 – Diagrama de Classes do Processo de Remoção de <i>Stopwords</i></b>	<b>107</b>
<b>Figura 8.4 – Diagrama de Seqüência do Processo de Remoção de <i>Stopwords</i></b>	<b>108</b>

<b>Figura 9.1 – Diagrama de Classes do Pacote <i>search</i></b>	<b>110</b>
<b>Figura 9.2 – Diagrama de Classes do Pacote <i>nlp</i></b>	<b>111</b>
<b>Figura 9.3 – Diagrama de Classes <i>CorpusSearch</i></b>	<b>112</b>
<b>Figura 9.4 – Diagrama de Classes <i>StemmingSearch</i></b>	<b>113</b>
<b>Figura 9.5 – Diagrama de Classes da Visão Geral do <i>Framework NLPSearch</i></b>	<b>114</b>
<b>Figura 10.1 – Aplicação <i>Stemming Search</i></b>	<b>116</b>
<b>Figura 10.2 – Resultado da Aplicação <i>Stemming Search</i></b>	<b>117</b>
<b>Figura 10.3 – Aplicação <i>Corpus Search</i></b>	<b>118</b>
<b>Figura 10.4 – Corpus gerado através da aplicação <i>Corpus Search</i></b>	<b>119</b>
<b>Figura 10.5 – Salvando o <i>Corpus</i> gerado</b>	<b>119</b>
<b>Figura 10.6 – <i>Corpus</i> salvo</b>	<b>120</b>
<b>Figura 10.7 – Estrutura de diretórios</b>	<b>121</b>
<b>Figura 10.8 – Estrutura do <i>Framework NLPSearch</i></b>	<b>122</b>
<b>Figura 10.9 – Estrutura de classes do pacote <i>search</i></b>	<b>122</b>
<b>Figura 10.10 - Estrutura de classes do pacote <i>nlp</i></b>	<b>123</b>
<b>Figura 10.11 - Estrutura de classes do pacote <i>nlpSearch</i></b>	<b>123</b>
<b>Figura 10.12 – Documentação do <i>Framework NLPSearch</i></b>	<b>124</b>



## LISTA DE TABELAS

<b>Tabela 2.1 - Comparativo dos Diagramas da <i>UML</i> 1.4 e 2.0 (MELO, 2004)</b>	<b>29</b>
<b>Tabela 2.2 – Diagramas divididos por categorias (MELO, 2004)</b>	<b>30</b>
<b>Tabela 4.1 – Parâmetros de Pesquisa do <i>Google</i> (GOOGLE, 2006)</b>	<b>55</b>
<b>Tabela 4.2 – Limitações da <i>API</i> do <i>Google</i> (GOOGLE, 2006)</b>	<b>55</b>
<b>Tabela 4.3 – <i>Tags</i> retornadas pelo <i>Google</i> (GOOGLE, 2006)</b>	<b>57</b>
<b>Tabela 4.4 – Parâmetros de Pesquisa do <i>Yahoo</i> (YAHOO, 2006)</b>	<b>59</b>
<b>Tabela 4.5 – <i>Tags</i> retornadas pelo <i>Yahoo</i> (YAHOO, 2006)</b>	<b>60</b>
<b>Tabela 4.6 – Parâmetros de Pesquisa do <i>Technorati</i> (TECHNORATI, 2006)</b>	<b>62</b>
<b>Tabela 4.7 – <i>Tags</i> retornadas pelo <i>Technorati</i> (TECHNORATI, 2006)</b>	<b>63</b>
<b>Tabela 4.8 – Comparativo entre os Mecanismos de Busca</b>	<b>63</b>

## LISTA DE ABREVIATURAS E SIGLAS

ADTF	<b>A</b> nalysys and <b>D</b> esign <b>T</b> ask <b>F</b> orce
API	<b>A</b> pplication <b>P</b> rogram <b>I</b> nterface
FAQ	<b>F</b> requently <b>A</b> sksed <b>Q</b> uestion
HTML	<b>H</b> iper <b>T</b> ext <b>M</b> arkup <b>L</b> anguage
IP	<b>I</b> nternet <b>P</b> rotocol
MIME	<b>M</b> ultipurpose <b>I</b> nternet <b>M</b> ail <b>E</b> xtensions
MVC	<b>M</b> odel <b>V</b> iew <b>C</b> ontroller
ODP	<b>O</b> pen <b>D</b> irectory <b>P</b> roject
OMG	<b>O</b> bject <b>M</b> anagement <b>G</b> roup
OMT	<b>O</b> bject <b>M</b> odeling <b>T</b> echnique
OOSE	<b>O</b> bject- <b>O</b> riented <b>S</b> oftware <b>E</b> ngineering
ORB	<b>O</b> bject <b>R</b> equest <b>B</b> roker
PDF	<b>P</b> ortable <b>D</b> ocument <b>F</b> ormat
PLN	<b>P</b> rocessamento da <b>L</b> inguagem <b>N</b> atural
RSS	<b>R</b> ich <b>S</b> ite <b>S</b> ummary/ <b>R</b> eally <b>S</b> imple <b>S</b> yndication
RTF	<b>R</b> evision <b>T</b> ask <b>F</b> orce
SGML	<b>S</b> tandard <b>G</b> eneralized <b>M</b> arkup <b>L</b> anguage
SMTP	<b>S</b> imple <b>M</b> ail <b>T</b> ransfer <b>P</b> rotocol
UML	<b>U</b> nified <b>M</b> odeling <b>L</b> anguage
UTF-8	<b>8</b> -bit <b>U</b> nicode <b>T</b> ransformation <b>F</b> ormat
XML	<b>E</b> Xtensible <b>M</b> arkup <b>L</b> anguage

## SUMÁRIO

<b>INTRODUÇÃO</b>	<b>14</b>
<b>1 ORIENTAÇÃO A OBJETOS</b>	<b>17</b>
1.1 Objetos	17
1.2 Mensagens	19
1.3 Classes	20
1.4 Abstração	21
1.5 Encapsulamento	22
1.6 Herança	23
1.7 Polimorfismo	25
<b>2 UML</b>	<b>27</b>
2.1 Sobre a <i>UML</i>	27
2.2 Diagramas da <i>UML</i>	28
2.2.1 Diagrama de Casos de Uso	30
2.2.2 Diagrama de Classes	31
2.2.3 Diagrama de Seqüência	32
2.2.4 Diagrama de Objetos	33
2.2.5 Diagrama de Atividades	34
2.2.6 Diagrama de Componentes	35
2.2.7 Diagrama de Implantação	35
2.2.8 Diagrama de Comunicação	36
2.2.9 Diagrama de Máquina de Estados	37
2.2.10 Diagrama de Pacotes	38
2.2.11 Diagrama de Estrutura Composta	38
2.2.12 Diagrama de Visão Geral	39
2.2.13 Diagrama Temporal	40
<b>3 FRAMEWORKS ORIENTADOS A OBJETOS</b>	<b>41</b>
3.1 <i>Framework</i> X Bibliotecas	42
3.2 <i>Framework Model View Controler (MVC)</i>	42
3.3 Classificação dos <i>Frameworks</i>	43
3.3.1 Quanto ao Propósito	44
3.3.2 Quanto as Técnicas de Reuso	45
3.4 Compreensão e Uso de <i>Frameworks</i>	46
3.5 Componentes Envolvidos no Desenvolvimento e Uso de <i>Frameworks</i>	47
3.6 Construção de <i>Frameworks</i>	49

3.7 Documentação de <i>Frameworks</i>	50
3.8 Vantagens e Desvantagens do Uso de <i>Frameworks</i>	51
<b>4 API'S DE MECANISMOS DE BUSCA</b>	<b>53</b>
4.1 <i>Google API</i>	53
4.2 <i>Yahoo API</i>	57
4.3 <i>Technorati API</i>	60
4.4 Comparativo entre as <i>API's</i> de Mecanismos de Busca	63
<b>5 PROCESSAMENTO DA LINGUAGEM NATURAL (PLN)</b>	<b>64</b>
5.1 Conceitos de PLN	64
5.2 Histórico do PLN	65
5.3 <i>Tagger, Morphological Analyser</i> (Etiquetador, Analisador Morfológico)	66
5.4 <i>Stemmer</i> (Radicalizador de Palavras)	67
5.5 <i>Parser</i> (Analisador Sintático)	68
5.6 <i>Corpus Tools</i>	69
5.6.1 Tipos de <i>Corpus</i>	70
5.6.2 Pré-Requisitos para Formação de um <i>Corpus</i> Computadorizado	70
<b>6 FERRAMENTAS DE PLN</b>	<b>72</b>
6.1 <i>WebJspell</i>	72
6.2 <i>Snowball</i>	75
6.3 <i>FLIP On-Line</i>	76
6.4 <i>JBootCat</i>	78
6.5 <i>API Apache Lucene</i>	83
<b>7 MODELAGEM DAS API'S DE MECANISMOS DE BUSCA</b>	<b>84</b>
7.1 Modelagem <i>Google API</i>	84
7.1.1 Realizar consultas através da <i>Google API</i>	84
7.1.2 Diagrama de Classes da <i>Google API</i>	85
7.1.3 Diagrama de Seqüência da <i>Google API</i>	86
7.2 Modelagem <i>Yahoo API</i>	86
7.2.1 Realizar consultas através da <i>Yahoo API</i>	87
7.2.2 Diagrama de Classes da <i>Yahoo API</i>	88
7.2.3 Diagrama de Seqüência da <i>Yahoo API</i>	89
7.3 Modelagem da <i>Technorati API</i>	89
7.3.1 Realizar consultas através da <i>Technorati API</i>	89
7.3.2 Diagrama de classes da <i>Technorati API</i>	90
7.3.3 Diagrama de Seqüência da <i>Technorati API</i>	91
<b>8 MODELAGEM DAS FERRAMENTAS DE PLN</b>	<b>92</b>
8.1 Modelagem de Casos de Uso da Ferramenta <i>WebJspell</i>	92
8.1.1 Informar a frase que será processada e alguns parâmetros iniciais	92
8.1.2 Análise morfológica da frase informada	93

8.1.3 Sugerir Análise	94
8.1.4 Remover do dicionário	94
8.1.5 Flexionar a palavra	95
8.1.6 Exemplos de frases	95
8.1.7 Ocorrências da palavra	96
8.1.8 Conjugação verbal	97
8.2 Modelagem de Casos de Uso da Ferramenta <i>Snowball</i>	97
8.2.1 Informar as palavras para realizar <i>stemming</i> ou lematização	97
8.2.2 Realizar <i>stemming</i> ou lematização	98
8.3 Modelagem de Casos de Uso da Ferramenta <i>Flip-Online</i>	98
8.3.1 Informar a frase que será processada e os parâmetros iniciais	98
8.3.2 Análise sintática da frase informada	99
8.4 Modelagem de Casos de Uso da Ferramenta <i>JBootCat</i>	100
8.4.1 Informar arquivo para armazenamento do <i>Corpus</i>	100
8.4.2 Informar termos de pesquisa	100
8.4.3 Montar tuplas de pesquisa	101
8.4.4 Realizar pesquisa na <i>web</i>	102
8.4.5 Selecionar resultados da pesquisa	102
8.4.6 Gerar corpus	103
8.5 Modelagem da API <i>Lucene</i>	104
8.5.1 Algoritmo de <i>Stemming</i> da API <i>Lucene</i>	104
8.5.1.1 Realizar <i>Stemming</i> em um texto	104
8.5.1.2 Diagrama de Classes do Processo de <i>Stemming</i>	105
8.5.1.3 Diagrama de Seqüência do processo de <i>Stemming</i>	106
8.5.2 Algoritmo de Remoção de <i>Stopwords</i>	106
8.5.2.1 Realizar processo de Remoção de <i>Stopwords</i> de um texto	106
8.5.2.2 Diagrama de Classes do Processo de Remoção de <i>Stopwords</i>	107
8.5.2.3 Diagrama de Seqüência do processo de Remoção de <i>Stopwords</i>	108
<b>9 MODELAGEM DO FRAMEWORK – DIAGRAMAS DE CLASSES</b>	<b>109</b>
9.1 Pacote <i>search</i>	109
9.2 Pacote <i>nlp</i>	110
9.3 Diagrama de Classes do Pacote <i>nlpSearch</i>	111
9.4 Diagrama de Classes da Visão Geral do <i>Framework NLPSearch</i>	113
<b>10 PROTOTIPAÇÃO</b>	<b>115</b>
10.1 Funcionalidades Implementadas	115
10.1.1 Protótipo <i>Stemming Search</i>	115
10.1.2 Protótipo <i>Corpus Search</i>	117
10.2 Como utilizar o <i>framework</i> para desenvolver aplicações	120
10.2.1 Estrutura necessária	120
10.2.2 Estrutura de pacotes do <i>Framework NLPSearch</i>	121
10.2.3 Estrutura de Classes do <i>Framework NLPSearch</i>	122
10.2.4 Documentação do <i>Framework NLPSearch</i>	124
<b>CONSIDERAÇÕES FINAIS</b>	<b>125</b>
<b>REFERÊNCIAS BIBLIOGRÁFICAS</b>	<b>127</b>

## INTRODUÇÃO

A investigação científica de métodos e técnicas para o Processamento da Linguagem Natural (PLN) têm se desenvolvido muito nas últimas décadas, e especialmente nos últimos anos, com esforços focados em diferentes fenômenos da linguagem. O processamento de anáforas, sumarização automática, o tratamento de diálogos, e a tradução automática são alguns exemplos do emprego de recursos computacionais e o conhecimento de lingüistas na construção de sistemas inteligentes que consigam tratar de forma individualizada os fenômenos da linguagem (GASPERIN et al., 2003; PARDO et. al, 2003; NUNES et al., 1999; RUSSEL; NORVIG, 2004).

O estudo das primeiras técnicas de PLN são datadas da década de 50, e eram focadas especialmente na tradução automática de textos. Segundo Nunes et al. (1999) o marco inicial das pesquisas de PLN foi a distribuição de 200 cópias de uma carta conhecida como *Weaver Memorandum*, escrita por Warren Weaver, exímio conhecedor de criptografia computacional. Nesta carta ele convidava a comunidade científica a desenvolver projetos sobre um novo campo, conhecido como tradução automática.

Nesse contexto de evolução das técnicas de PLN surge a necessidade de desenvolver mecanismos para integração desses esforços, aliados a uma larga base de dados como a internet. Para propor um sistema que integre estes conceitos de forma a ampliar sua utilização em diferentes contextos lingüísticos, este trabalho propõe a construção de um *framework* para o desenvolvimento de aplicações de PLN que utilizem a *Web* como *corpus*<sup>1</sup>.

Existem diversos conceitos para definir um *framework* orientado a objetos, segundo Silva (2000) um *framework* orientado a objetos utiliza o paradigma de orientação a objetos para produzir uma descrição de um domínio para ser reutilizada. Sendo assim um *framework* é uma estrutura de classes inter-relacionadas, que correspondem a uma implementação

---

<sup>1</sup> Corpus – “Grande coleção de texto, como bilhões de páginas que constituem a web” (RUSSELL; NORVIG, 2004, p. 807).

incompleta para um conjunto de aplicações de um domínio, sendo que esta estrutura de classes deve ser adaptada para a geração de aplicações específicas.

Utilizar um *framework* consiste em criar aplicações a partir deste *framework*. A motivação para se usar um *framework* na construção de aplicações é a perspectiva de um aumento considerável da produtividade e da qualidade do sistema desenvolvido, em função da reutilização promovida. Para produzir uma aplicação utilizando um *framework* como base, deve-se estender sua estrutura de modo a cumprir os requisitos básicos da aplicação a que se quer desenvolver (SILVA, 2000).

Um dos obstáculos encontrados para se desenvolver aplicações através do uso de *frameworks* é justamente compreender seu uso, como funciona e interage internamente. A questão é que se um desenvolvedor de aplicações necessitar empreender um esforço considerável para aprender a utilizar um *framework*, que para desenvolver uma aplicação do início, a utilidade deste *framework* acaba se tornando questionável (SILVA, 2000).

Para montar um *corpus* é necessário à reunião de diversos textos, para isto este trabalho propõe que este corpus seja montado utilizando a base de dados contida na *web*. Para se ter acesso a essa base de dados alguns mecanismos de busca disponibilizam gratuitamente suas *API's*<sup>2</sup> de pesquisa na *web* para desenvolvimento de aplicações não comerciais. O presente trabalho irá utilizar as *API's* disponibilizadas pelo *Google*, *Yahoo* e *Technorati* para realizar pesquisas e montar o *corpus* necessário para o processamento das técnicas de PLN.

Dessa forma a principal proposta deste trabalho é a construção de um *framework* orientado a objetos para o desenvolvimento de aplicações para PLN que utilizem a *web* como *corpus*, através da utilização de *API's* de mecanismos de busca conhecidos.

O presente trabalho está dividido em dez capítulos, organizados seguinte forma:

O primeiro capítulo mostra uma introdução sobre a utilização da orientação a objetos para construção de *softwares*, este capítulo torna-se muito interessante, pois serve como base para o desenvolvimento do *framework* proposto. O segundo capítulo aborda a utilização da

---

<sup>2</sup> API – Do Inglês *Application Programming Interface* (*Interface* de Programação de Aplicativos). É um conjunto de rotinas e padrões estabelecidos por um software para utilização de suas funcionalidades por programas aplicativos, isto é, programas que não querem envolver-se em detalhes da implementação do *software*, mas apenas usar seus serviços.

*UML* para a modelagem de sistemas de *software* orientados a objetos. O terceiro capítulo apresenta uma abordagem sobre *frameworks* orientados a objetos, onde são apresentadas vantagens e desvantagens da sua utilização.

O quarto capítulo apresenta algumas noções dos mecanismos de busca que são utilizados no trabalho.

No quinto capítulo são apresentados alguns conceitos sobre PLN, assim como algumas definições a respeito das suas principais ferramentas. No sexto capítulo são apresentadas algumas ferramentas que se dispõem a atender os conceitos abordados no quinto capítulo.

No sétimo e oitavo capítulos deste trabalho são apresentadas as propostas de análise e modelagem das *API's* dos mecanismos de busca anteriormente citadas (quarto capítulo), assim como das ferramentas de PLN discutidas (sexto capítulo). Estes dois capítulos são importantes principalmente por introduzirem algumas noções sobre conceitos que estão presentes no *Framework NLPSearch*.

No nono capítulo é apresentada a modelagem do *framework* baseada unicamente em diagramas de classes. A modelagem está baseada somente neste tipo de diagrama porque nos capítulos anteriores foram modeladas as principais funcionalidades que estão presentes no produto final deste trabalho.

O décimo capítulo está dividido em duas principais seções. Na primeira seção são apresentados os protótipos desenvolvidos através da utilização da estrutura do *Framework NLPSearch*. Na segunda seção primeiramente são apresentados os principais requisitos para se utilizar o *framework* como base para a construção de novas aplicações. Em seguida são apresentadas as estruturas que fazem parte da arquitetura do *framework*.



# 1 ORIENTAÇÃO A OBJETOS

O desenvolvimento de um sistema de *software* não é tarefa simples de ser realizada, pois demanda muita lógica computacional, conhecimento de linguagens de programação, ferramentas de desenvolvimento e por último, mas não com menos importância, do negócio a que se destina o sistema a ser desenvolvido. O programador ou equipe responsável pelo desenvolvimento precisam ter esse conhecimento para produzir *softwares* que primem essencialmente pela qualidade. Hoje em dia para atingir essa qualidade tornou-se essencial a utilização do paradigma<sup>3</sup> de orientação a objetos no desenvolvimento de sistemas de *software*. As técnicas orientadas a objetos prometem aumentar de forma significativa tanto a qualidade quanto a produtividade do desenvolvimento de *software* (YOURDON; ARGILA, 1999). Segundo Bezerra (2002) um sistema de *software* orientado a objetos consiste de objetos em colaboração com o objetivo de realizar as funcionalidades desse sistema. Cada objeto é responsável por tarefas específicas. É através da cooperação entre objetos que a computação do sistema se desenvolve. Neste capítulo são apresentadas algumas noções sobre o paradigma de orientação a objetos.

## 1.1 Objetos

O paradigma da orientação a objetos visualiza um sistema de *software* como uma coleção de agentes interconectados chamados objetos. Cada objeto é responsável por realizar tarefas específicas. É através da interação entre objetos que uma tarefa computacional é realizada (BEZERRA, 2002).

---

<sup>3</sup> “Paradigma. [Do grego parádeigma, pelo latim tard. paradgma.] S. m. Termo com o qual Thomas Kuhn designou as realizações científicas (p. ex., a dinâmica de Newton ou a química de Lavoisier) que geram modelos que, por período mais ou menos longo e de modo mais ou menos explícito, orientam o desenvolvimento posterior das pesquisas exclusivamente na busca da solução para os problemas por elas suscitados” (BEZERRA, 2002, p. 4).

Paradigma é Definido como um “conjunto de teorias, padrões e métodos que juntos representam um modo de organizar conhecimento” (LEE; TEPFENHART, 2001, p. 24).

“Objeto é uma entidade independente, assíncrona e concorrente que ‘sabe coisas’ (armazena dados), ‘realiza trabalho’ (encapsula serviços) e ‘colabora com outros objetos’ (por meio de troca de mensagens) para executar as funções finais do sistema” (YOURDON; ARGILA, 1999, p. 7).

Segundo Harmon e Watson (1998) um objeto é um complexo tipo de dado que contém outros tipos de dados, chamados de atributos ou variáveis, e módulos de código, chamados operações ou métodos. Os atributos e valores associados estão escondidos dentro do objeto. Qualquer outro objeto que quiser obter ou alterar um valor associado com um primeiro objeto deve fazer isso enviando uma mensagem para ele.

Pree (1995) define um objeto como sendo uma instância de tipo de dados abstratos. Desta maneira um objeto é uma entidade que contém atributos (dados representando um estado do objeto) e que disponibiliza certas operações que são definidas para um objeto em particular. Os atributos são chamados de variáveis de instância e as operações são chamadas de métodos.

Uma definição mais simples de objeto é apresentada por Rumbaugh et al. (1994) que diz que um objeto é simplesmente alguma coisa que faz sentido no contexto de uma aplicação e ainda define um objeto como sendo um conceito, uma abstração, algo com limites nítidos e significado em relação ao problema em causa. Objetos servem a dois objetivos: eles facilitam a compreensão do mundo real e oferecem uma base real para a implementação em um computador.

Para Coleman (1996) um objeto representa simplesmente um elemento que pode ser identificado de maneira única. Em um nível apropriado de abstração praticamente tudo pode ser considerado como sendo um objeto. Dessa forma, elementos específicos como pessoas, empresas, acessórios, automóveis ou eventos podem ser considerados como objetos.

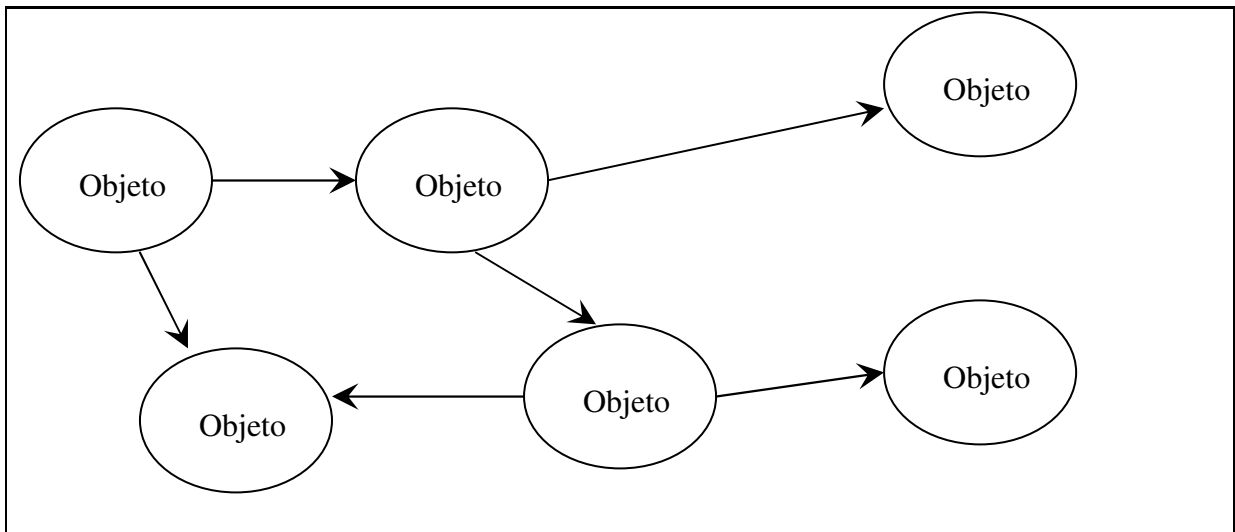
Objetos possuem identidades e são distinguíveis. O termo identidade significa que os objetos se distinguem um dos outros pela sua própria existência, ou seja, mesmo possuindo os mesmos atributos são considerados distintos, pois cada qual representa uma instância diferente de uma mesma classe de objetos (RUMBAUGH et al., 1994).

Quando um objeto quer mudar ou solicitar alguma informação de outro objeto ele envia mensagens, estas mensagens servem como gatilhos para disparar métodos para o outro objeto, são exatamente estes métodos que irão realizar a tarefa solicitada pelo primeiro objeto. Os métodos que podem ser chamados por outros objetos são conhecidos como métodos públicos. Objetos ainda podem conter métodos privados (métodos visíveis somente para o próprio objeto), não sendo possível acessá-los através do envio de mensagens (LEE; TEPFENHART, 2001).

## **1.2 Mensagens**

Objetos não executam suas operações (métodos) aleatoriamente. Para que um método seja executado deve haver um estímulo enviado a este objeto. Um objeto é visto como uma abstração do mundo real, dessa forma faz sentido dizer que um objeto pode responder a estímulos a ele enviados. Independentemente da origem do estímulo, quando ele ocorre, diz-se que o objeto em questão está recebendo uma mensagem. Esta mensagem é recebida pelo objeto e este objeto deve realizar alguma operação (BEZERRA, 2002). Quando uma mensagem é enviada a um objeto, o método específico que será executado depende de ambos, ou seja, do objeto que enviou a mensagem e do objeto receptor. A associação em tempo de execução de uma solicitação de um objeto e uma das suas operações é conhecida como ligação dinâmica (GAMMA et al., 2000). Para simplificar a compreensão desse processo de troca de mensagens entre objetos pode-se considerar um esquema cliente-servidor, onde um objeto cliente solicita alguma informação ao objeto servidor, esse por sua vez processa a solicitação e devolve o resultado para o objeto cliente.

A figura 1.1 representa a troca de mensagens entre objetos.



**Figura 1.1 - Objetos interagem através do envio de mensagens (BEZERRA, 2002)**

### 1.3 Classes

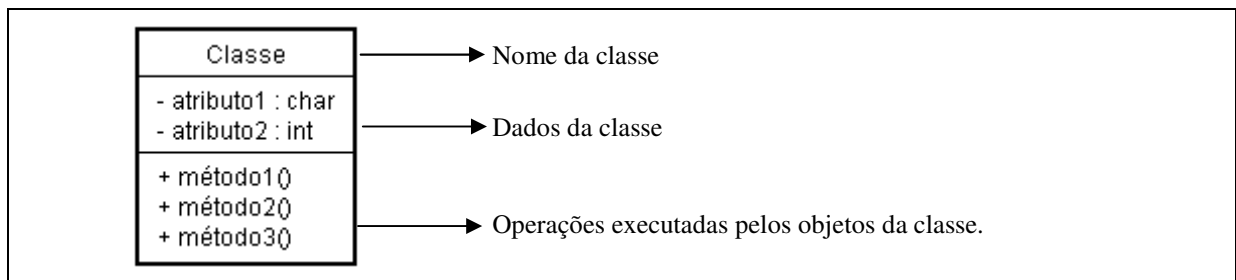
Uma classe de objetos descreve um grupo de objetos com propriedades semelhantes (atributos), o mesmo comportamento (métodos), os mesmos relacionamentos com outros objetos e a mesma semântica (RUMBAUGH et al., 1994).

Harmon e Watson (1998) definem uma classe como sendo um modelo que possui métodos e atributos e tipos de informações, mas não carregam valores atuais. Os objetos gerados pelas classes carregam valores atuais. Algumas classes são muito abstratas e são simplesmente utilizadas para prover uma estrutura para outras classes, entretanto grande parte das classes são modelos usados para gerar os objetos que serão responsáveis pela execução do sistema.

As classes são blocos de construção da maior parte das linguagens orientadas a objetos (COLEMAN et al., 1996).

A principal finalidade das classes abstratas é definir uma *interface* para suas subclasses. Um classe abstrata postergará parte de, ou toda, sua implementação para os métodos definidos em suas subclasses, sendo assim uma classe abstrata não pode ser instanciada. Os métodos que uma classe abstrata declara, mas não implementa, são chamados de métodos abstratos. As classes que não são abstratas são conhecidas como classes concretas. Estas classes ao contrário das abstratas podem ser instanciadas (GAMMA et al., 2000).

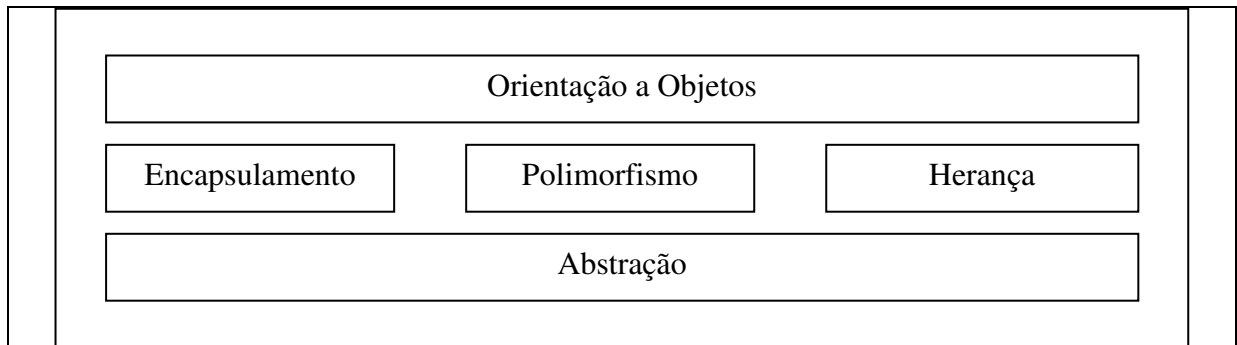
Novas classes podem ser instanciadas a partir de classes existentes através de um dos princípios mais importantes da orientação a objetos, conhecido como herança de classes. As subclasses instanciadas a partir de uma classe mãe herdam todos os atributos e métodos que a classe mãe define. Os objetos que são instâncias das subclasses irão conter todos os dados definidos pela subclasse e suas classes mãe, e eles serão capazes de executar todas as operações definidas por esta subclasse e seus “ancestrais” (GAMMA et al., 2000). A notação gráfica de uma classe é representada como sendo uma “caixa” com no máximo três compartimentos exibidos, dependendo do grau de abstração. No primeiro é exibido o nome da classe, no segundo aparecem os atributos e no terceiro as operações ou métodos, conforme mostra a figura 1.2 (BEZERRA, 2002).



**Figura 1.2 - Representação Gráfica de Uma Classe (BEZERRA, 2002)**

#### 1.4 Abstração

A abstração consiste em concentra-se nos aspectos essenciais, próprios, de uma entidade e ignorar as propriedades acidentais. Abstrações permitem gerenciar a complexidade e concentrar a atenção nas características essenciais dos objetos. O uso da abstração preserva a liberdade de se tomar decisões evitando, tanto quanto possível, comprometimentos prematuros com detalhes. Uma abstração é dependente da perspectiva, ou seja, o que pode ser considerado importante em um determinado contexto pode não ter tal importância em outro (RUMBAUGH et al., 1994; BEZERRA, 2002). Levando em consideração essa definição, pode-se afirmar que o Princípio da Abstração é aplicado em todos os demais Princípios da Orientação a Objetos, como mostra a figura 1.3.



**Figura 1.3 - Princípios da orientação a objetos como aplicações do Princípio da Abstração (BEZERRA, 2002)**

### 1.5 Encapsulamento

Também chamado de “ocultar informações”, este conceito estabelece que cada objeto contenha todos os detalhes de implementação necessários sobre como ele funciona e oculta todos os detalhes sobre como ele executa suas operações ou gerencia sua estrutura interna (REED, 2000).

Segundo Bezerra (2002) o mecanismo de encapsulamento é uma forma de restringir o acesso ao comportamento interno de um objeto. Um objeto que precise trocar informações com outro objeto para desempenhar alguma tarefa simplesmente envia uma mensagem para este último. O método requisitado não é conhecido dos objetos que requisitam alguma informação.

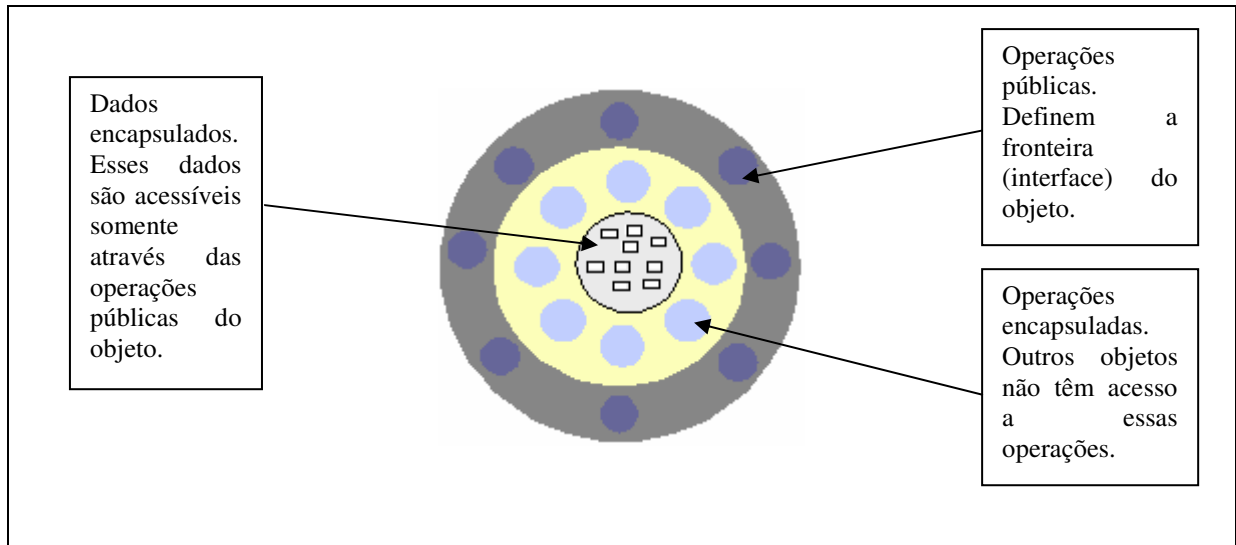
O encapsulamento impede que um programa se torne tão interdependente que uma simples modificação possa causar grandes efeitos de propagação, dessa forma a implementação de um objeto pode ser modificada sem que essa modificação afete as aplicações que o utilizem (RUMBAUGH et al., 1994).

Para que um objeto requisitante saiba quais tarefas que um outro objeto sabe realizar ou quais informações ele conhece, este precisa ter conhecimento da *interface* deste outro objeto. A *interface* de um objeto define os serviços que ele sabe realizar e conseqüentemente as mensagens que ele recebe quando é solicitado (BEZERRA, 2002).

Sendo assim, o encapsulamento também pode ser visto e “comparado” com o conceito de caixa-preta – Um bloco de informações desconhecidas. Simplesmente, utiliza-se este bloco, fornecendo entradas e recebendo saídas, sem saber como essas entradas são

processadas para fornecer as saídas. Não interessa o que tem dentro do bloco. O que interessa é que este bloco fornece as saídas desejadas, de acordo com as entradas fornecidas.

A figura 1.4 ilustra um objeto. Os dados no centro podem ser usados apenas com as operações constantes do anel exterior.



**Figura 1.4 - Objeto – Conceito de Encapsulamento**

## 1.6 Herança

A Herança pode ser considerada como uma abstração poderosa para o compartilhamento de similitudes entre classes, ao mesmo tempo em que suas diferenças são preservadas (RUMBAUGH et al., 1994). No conceito de herança as classes semelhantes são agrupadas em hierarquias. Cada nível hierárquico pode ser visto como um nível de abstração, sendo assim, cada classe em um nível de hierarquia herda as características das classes que estão em um nível mais alto (BEZERRA, 2002). A utilização da herança possibilita aos desenvolvedores alavancar o poder da hierarquia de classes criando novos objetos sem ter que desenvolvê-los novamente do zero (REED, 2000).

As classes que recebem a herança são denominadas subclasses, enquanto as classes que geram as subclasses são denominadas superclasses. Os atributos e operações comuns a um agrupamento de subclasses são incluídos na superclasse e são compartilhados por todas as subclasses, por isso diz-se que cada subclasse herda todas as características da sua superclasse (RUMBAUGH et al., 1994).

[...] Uma instância de uma subclasse é simultaneamente uma instância de todas as suas classes ancestrais. O estado de uma instância inclui um valor para cada atributo de cada classe ancestral. Qualquer operação em qualquer classe ancestral pode ser aplicada a uma instância. Cada subclasse não só herda todas as características de seus ancestrais como também acrescenta seus próprios atributos específicos e suas próprias operações [...] (RUMBAUGH et al., 1994, p. 54).

Existem dois tipos de heranças, a herança de classe (implementação) e a herança de *interface* (subtipificação). Segundo Gamma et al. (2000) é importante compreender a diferença entre estes dois tipos. A herança de classe define a implementação de um objeto em termos da implementação de outro, ou seja, pode ser considerado como um mecanismo para compartilhar código e representação. A herança de *interface* descreve quando um objeto pode substituir outro.

A figura 1.5 mostra um exemplo de herança para figuras gráficas, onde as classes em um nível hierárquico inferior herdam da classe que está em um nível hierárquico superior. Analisando o diagrama abaixo é possível verificar como funciona a herança de operações: Mover, selecionar, girar e exibir são operações herdadas por todas as subclasses. Medir aplica-se a figuras uni e bidimensionais. Preencher aplica-se somente a figuras bidimensionais (RUMBAUGH, et al., 1994).

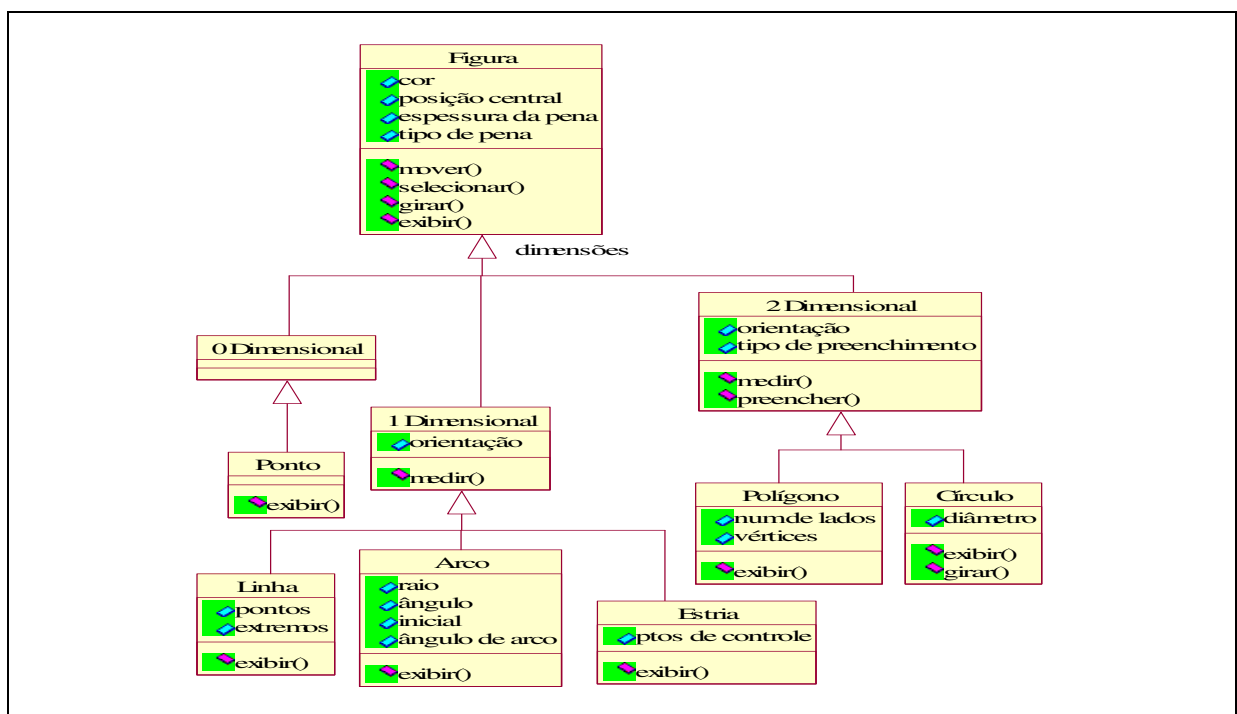


Figura 1.5 - Herança para Figuras Gráficas (RUMBAUGH, et al., 1994)



## 1.7 Polimorfismo

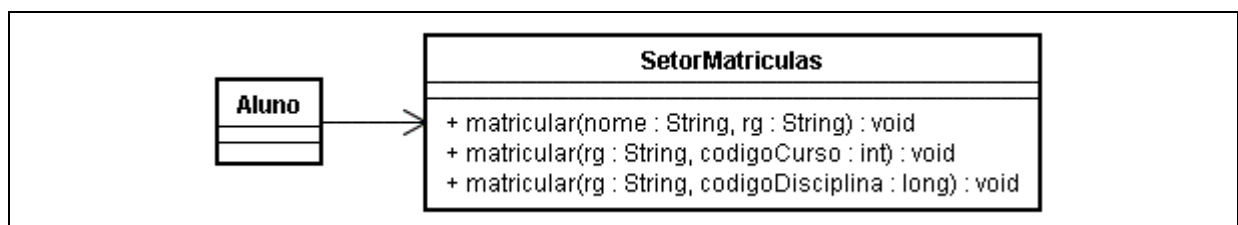
A palavra polimorfismo é originada do grego, onde “poli” significa muitas e “morfismo” significa formas, então polimorfismo significa muitas formas. Em se tratando de orientação a objetos polimorfismo é conhecido como um importante conceito, pois através dele uma mesma operação/método pode atuar de modos diversos em classes distintas (RUMBAUGH et al., 1994).

Segundo Bezerra (2002) o polimorfismo indica a capacidade de abstrair várias implementações diferentes em uma única *interface*.

Para Harmon e Watson (1998) polimorfismo significa que operações iguais poderão se comportar de forma diferente quando aplicadas a objetos de diferentes classes. Isso também quer dizer que diferentes operações associadas com diferentes classes podem interpretar uma mesma mensagem de diversas formas.

Reed (2000) também reforça a idéia de polimorfismo como sendo um importante conceito na orientação a objetos, pois segundo ele é importante poder enviar a mesma mensagem para diferentes objetos e fazê-los responder de forma correta.

A figura 1.6 apresenta um exemplo de polimorfismo aplicado a uma classe chamada de “SetorMatriculas”. Neste exemplo um aluno, representado pela classe “Aluno”, pode solicitar a realização de sua matrícula em um determinado curso ou em determinadas disciplinas para o Setor de Matriculas. Contudo, a forma como a matrícula será realizada vai depender do nível de ensino do curso que o aluno quer se matricular: extensão, especialização, graduação, etc.



**Figura 1.6 - Exemplo de Polimorfismo**

Para o aluno realizar a matrícula em um curso de extensão será necessário o seu nome e seu RG. Para uma matrícula num curso de especialização é necessário o RG e a identificação do curso. E por fim, para realizar a matrícula num curso de graduação será

necessário o RG e o código da disciplina que o aluno deseja cursar. Somente através do polimorfismo é possível realizar a mesma ação, no caso do exemplo, realizar a ação de “Matricular” de várias formas distintas, dando maior flexibilidade ao sistema.

Saindo um pouco do conceito teórico e entrando em questões de programação, percebe-se que a implementação do polimorfismo é a construção de métodos de classe com o mesmo nome (mesma ação), porém com assinaturas diferentes (mesma ação de forma variada). O conjunto de argumentos e tipos de retorno é chamado de assinatura (REED, 2000). No exemplo citado se vê três métodos com nomes iguais, mas com assinaturas distintas.

O segundo capítulo traz algumas noções sobre a utilização da Linguagem Unificada de Modelagem (*UML*), para modelagem de sistemas orientados a objetos.

## 2 UML

Segundo Booch et al. (2000) a *UML* é uma linguagem padrão para a elaboração da estrutura de projetos de *software* e pode ser empregada para a visualização, especificação, construção e a documentação de artefatos que façam uso de sistemas complexos de *software*. A *UML* é considerada uma linguagem de modelagem visual, ou seja, é um conjunto de notações e semântica correspondente com o intuito de representar visualmente uma ou mais perspectivas de um sistema (BEZERRA, 2002).

A *UML* é adequada para a modelagem de sistemas, cuja abrangência poderá incluir sistemas de informação corporativos a serem distribuídos a aplicações baseadas em web e até sistemas complexos embutidos de tempo real. É uma linguagem muito expressiva, abrangendo todas as visões necessárias ao desenvolvimento e implantação desses sistemas [...] (BOOCH et al., 2000, p. 13).

### 2.1 Sobre a UML

Com o passar dos anos a modelagem de sistemas passou a ser um ponto fundamental na estruturação de projetos bem sucedidos. Por este motivo iniciou-se uma busca por padrões de desenvolvimento de *software* que agregassem qualidade ao produto desenvolvido. Durante o fim da década de 80 e início da década de 90 surgiram diversas propostas de técnicas para modelagem de sistemas orientados a objetos. No período de 1989 a 1994 a quantidade de métodos orientados a objetos aumentou de um pouco mais de 10 para mais de 50 (BOOCH et al., 2000). No entanto a comunidade de *software* necessitava de um padrão que fosse aceito tanto pelas indústrias como pelo meio acadêmico.

Sendo assim, em 1996 surge a Linguagem de Modelagem Unificada para justamente unificar e padronizar as técnicas de modelagem existentes. Para a construção da *UML* foram estudadas várias notações, porém a base fundamental foi retirada dos estudos de três pesquisadores, Grady Booch, James Rumbaugh e Ivar Jacobson, frequentemente chamados de “os três amigos”. No processo de definição da *UML*, procurou-se aproveitar o melhor das características das notações preexistentes, principalmente das técnicas propostas

anteriormente pelos três amigos (essas técnicas eram conhecidas pelos nomes *Booch Method*<sup>4</sup>, *OMT*<sup>5</sup> e *OOSE*<sup>6</sup>) (BEZERRA, 2002).

Em janeiro de 1997 foi oferecida a *OMG*<sup>7</sup> (*Object Management Group*) a versão 1.0 da *UML* em resposta a solicitação da própria *OMG* para uma linguagem padrão de modelagem. Após a criação de um grupo de tarefas em semântica liderado por Cris Kobryn da *MCI Systemhouse* e administrado por Ed Eykholt da *Rational* com o propósito de formalizar a especificação da *UML* e de integrar a linguagem a outros esforços de padronizações a versão 1.1 da *UML* foi oferecida para padronização ao *OMG* em julho de 1997. Em setembro do mesmo ano essa versão foi aceita pela *ADTF* (*Analysys and Design Task Force*) e pelo *Architecture Board* do *OMG*, e foi submetido ao voto de todos os membros do *OMG*. A *UML* versão 1.1 foi adotada pelo *OMG* em 14 de novembro de 1997 como linguagem padrão (BOOCH et al., 2000; FOWLER; SCOTT, 2000; REED, 2000; BEZERRA, 2002). A manutenção da *UML* passou a ser responsabilidade da *RTF* (*Revision Task Force*) que pertence ao *OMG* e dirigida por Cris Kobryn. Dentre os objetivos desta força tarefa estão inseridas as revisões nas especificações da *UML* para propor novas alterações de forma que não provoquem mudanças no escopo original da linguagem. Em julho de 1998 foi apresentada a versão 1.2 da *UML* e no final do ano a versão 1.3. Em maio de 2001 foi lançada a versão 1.4, em agosto do mesmo ano foi submetido um relatório provisório sobre a versão 1.5, publicada em março de 2003. Entre março e junho de 2003 as primeiras especificações da *UML* 2.0 foram adotadas pelo *OMG* como padrão (MELO, 2004). A documentação oficial sobre a *UML* pode ser encontrada no site da *OMG*, [www.omg.org](http://www.omg.org).

## 2.2 Diagramas da *UML*

Um processo de desenvolvimento que utilize a *UML* como linguagem de suporte a modelagem envolve a criação de diversos documentos. Estes documentos podem ser textuais ou gráficos. Na terminologia da *UML*, estes documentos são conhecidos como artefatos de

---

<sup>4</sup> *Booch Method* – Expressivo especialmente durante as fases de projeto e construção de sistemas (BOOCH, 2000).

<sup>5</sup> *OMT* (*Object Modeling Technique*) – Mais útil para a análise e sistemas de informações com uso intensivo de dados (BOOCH, 2000).

<sup>6</sup> *OOSE* (*Object-Oriented Software Engineering*) – Fornecia excelente suporte para os casos como uma maneira de controlar a captura de requisitos, a análise e o projeto em alto nível (BOOCH, 2000).

<sup>7</sup> “*OMG* (*Object Management Group*) – Consórcio internacional de empresas que define e ratifica os padrões na área de orientação a objetos. ([www.omg.org](http://www.omg.org))” (BEZERRA, 2002, p. 13).

*software*, e são estes artefatos que compõe as diversas visões do sistema (BEZERRA, 2002). Os artefatos gráficos produzidos através da utilização da *UML* são chamados de diagramas.

Segundo Boch et al. (2000) um diagrama é uma representação gráfica de um conjunto de elementos, comumente representada como um gráfico conectado de vértices (itens) e arcos (relacionamentos). São desenhados para permitir a visualização de um sistema sob diferentes perspectivas, sendo assim um diagrama é considerado uma projeção em um sistema. Até a versão 1.4 a *UML* incluía nove diagramas em sua notação. A partir da versão 2.0 este número aumentou consideravelmente, sendo incluídos mais cinco diagramas, e alterado o nome de alguns existentes (MELO, 2004). A tabela 2.1 demonstra quais diagramas da *UML* 2.0 são novos e quais correspondem a diagramas da versão 1.4.

Diagramas da <i>UML</i> 1.4	Diagramas da <i>UML</i> 2.0
...	Pacotes
...	Estrutura Composta
...	Visão Geral
...	Temporal
Atividades	Atividades
Casos de Uso	Casos de Uso
Classes	Classes
Colaboração	Comunicação
Componentes	Componentes
Gráfico de Estado	Máquina de Estados
Implantação	Implantação
Objetos	Objetos
Seqüência	Seqüência

**Tabela 2.1 - Comparativo dos Diagramas da *UML* 1.4 e 2.0 (MELO, 2004)**

Os diagramas podem ser divididos em duas categorias, diagramas estáticos ou estruturais e diagramas dinâmicos ou comportamentais. A função dos diagramas estruturais é a de mostrar as características de um sistema que não mudam com o tempo, ou seja, servem para visualizar, especificar, construir e documentar os aspectos estáticos de um sistema (BOOCH, 2000; MELO, 2004). A função dos diagramas dinâmicos é mostrar as interações ativas que um sistema suporta, desta forma eles detalham a interação entre os diagramas estruturais (REED, 2000). A tabela 2.2 mostra os diagramas existentes na *UML* divididos nas duas categorias anteriormente citadas.

<b>Diagramas Estruturais</b>	Diagrama de Classes Diagrama de Objetos
------------------------------	--

	Diagrama de Componentes Diagrama de Pacotes Diagrama de Implantação Diagrama de Estrutura Composta
<b>Diagramas Dinâmicos</b>	Diagrama de Casos de Uso Diagrama de Interação: - Diagrama de Visão Geral - Diagrama de Seqüências - Diagrama Temporal - Diagrama de Comunicação Diagrama de Atividades Diagrama de Máquina de Estados

**Tabela 2.2 – Diagramas divididos por categorias (MELO, 2004)**

A seguir será apresentado um pequeno resumo sobre cada um dos diagramas que compõe a notação atual da *UML*.

### **2.2.1 Diagrama de Casos de Uso**

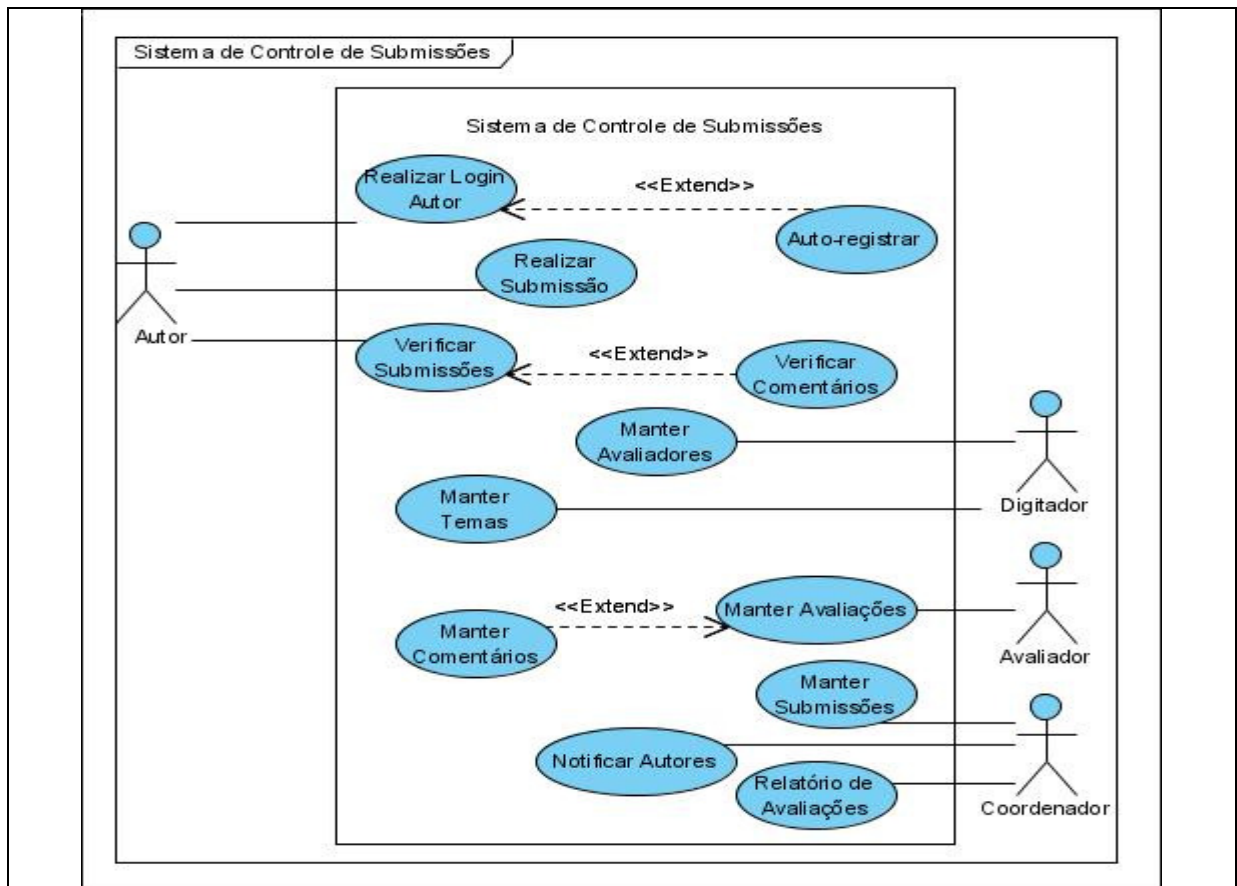
Quando utilizamos a *UML* para modelar um sistema, um dos primeiros processos a ser realizado diz respeito à modelagem de casos de uso. Um diagrama de caso de uso representa a interação entre usuário e sistema, demonstrando alguma funcionalidade relativa a eles (QUATRINI, 2001). Os diagramas de casos de uso são importantes para visualizar, especificar e documentar o comportamento de um elemento dentro do sistema (BOOCH et al., 2000).

Internamente, um caso de uso é uma seqüência de ações que permeiam a execução completa de um comportamento esperado para o sistema. Este é o ponto de partida para todo o processo, por isso é necessário ter muito cuidado e não incluir detalhes demais. Um Diagrama de Caso de Uso deve ser simples, porém não pode ser incompleto (BOOCH et al., 2000; BEZERRA, 2002; HARMON; WATSON, 1998; MELO, 2004). Desta forma um caso de uso não deve se preocupar com o “como” das funções. Este diagrama não deve prescrever seqüências de execução, mais sim as ações que são executadas pelos Atores.

Um ator está relacionado a um papel. No contexto do sistema, papel diz respeito à visão dada ao sistema (FOWLER, 2000; BEZERRA, 2002). Por exemplo, em um projeto de desenvolvimento de sistemas, um membro da equipe de desenvolvimento pode exercer vários papéis distintos, como projetista de testes, projetista de banco de dados, revisor de código, etc. Cada papel destaca funções diferentes específicas a ele. Isso deixa claro que a figura do ator

não está relacionada de forma direta a uma pessoa. Segundo Fowler (2000) um ator não precisa ser necessariamente humano. Um ator pode ser também um sistema externo que precisa de informações geradas em um sistema corrente.

A notação utilizada para ilustrar os atores de um caso de uso é a figura de um boneco, com o nome do ator definido abaixo dessa figura (a partir da versão 2.0 da *UML* é obrigatório informar o nome do ator). Cada caso de uso é representado por uma elipse, o nome do caso de uso é posicionado abaixo ou dentro da elipse. Um relacionamento de comunicação é representado por um segmento de reta ligando ator e caso de uso (BEZERRA, 2002; MELO, 2004). A figura 2.1 ilustra um Diagrama de Casos de Uso.



**Figura 2.1 – Diagrama de Casos de Uso (GUEDES, 2005)**

### 2.2.2 Diagrama de Classes

Um diagrama de classes mostra um conjunto de classes, *interfaces* e colaborações e seus relacionamentos. Os diagramas de classes são os diagramas mais encontrados em sistemas de modelagem orientados a objetos (BOOCH et al., 2000). Segundo Bezerra (2002),

o diagrama de classes é utilizado desde o nível de análise até o nível de especificação do projeto. De todos os diagramas da *UML*, esse é o mais rico em termos de notação.

O diagrama de classes demonstra a estrutura estática do *software* a ser construído, isto é, as classes, os relacionamentos (associações) entre suas instâncias (objetos), restrições e hierarquias. O diagrama é considerado estático, pois a estrutura descrita é sempre válida em qualquer ponto do ciclo da vida do sistema modelado (FALBO, 2001). A figura 2.2 mostra um exemplo de diagrama de classes.

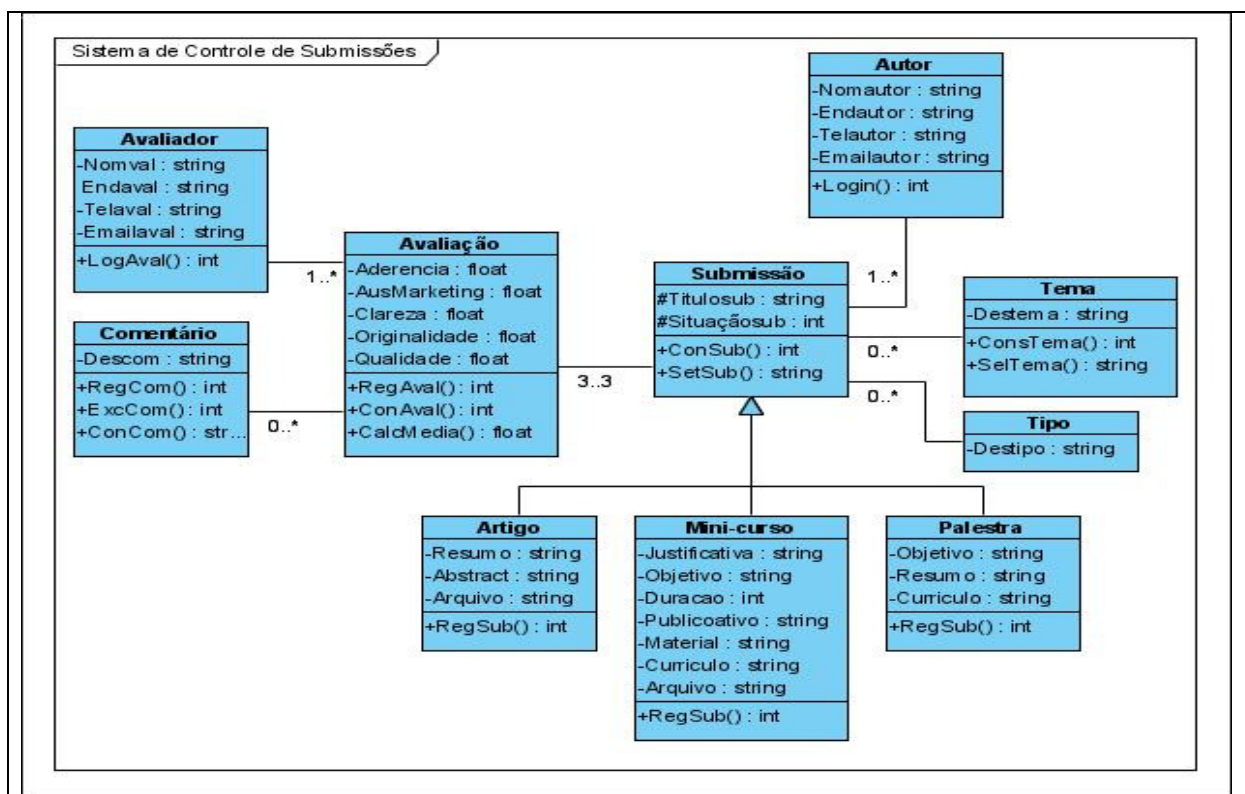


Figura 2.2 – Diagrama de Classes (GUEDES, 2005)

### 2.2.3 Diagrama de Seqüência

Segundo Booch et al. (2000) um diagrama de seqüência é um diagrama de interação que dá ênfase à ordenação temporal de mensagens. Um diagrama de seqüência mostra um conjunto de objetos e as mensagens enviadas e recebidas por estes objetos.

Os diagramas de seqüência são compostos por dois eixos: no eixo vertical é exibido o tempo e no eixo horizontal os objetos envolvidos na seqüência de certa atividade, como, por exemplo, na realização de um Cenário de Caso de Uso. Cada objeto é representado por um



retângulo e uma linha vertical pontilhada chamada de linha de vida do objeto, indicando a execução do objeto durante a seqüência. A comunicação entre os objetos é representada com uma linha com setas horizontais simbolizando as mensagens (BEZERRA, 2002). A figura 2.3 mostra um exemplo de digrama de seqüência.

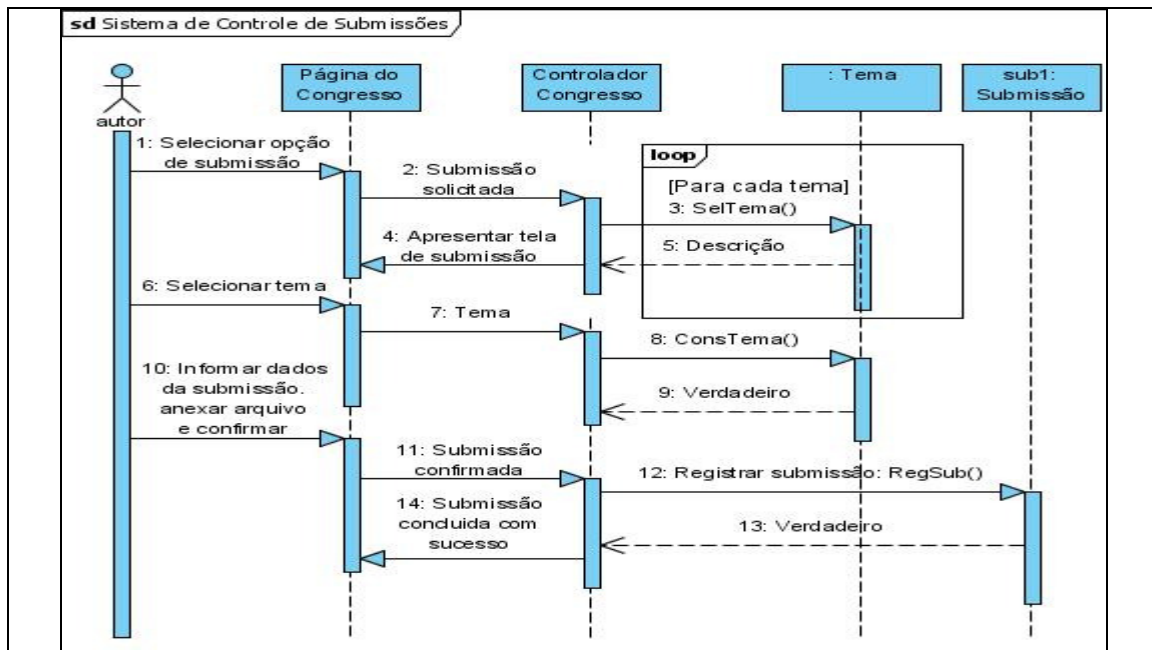


Figura 2.3 – Diagrama de Seqüência (GUEDES, 2005)

#### 2.2.4 Diagrama de Objetos

O diagrama de objetos mostra um conjunto de objetos e seus relacionamentos. Os diagramas de objetos direcionam a visão estática do projeto de um sistema ou a visão estática do processo de um sistema, tal qual os diagramas de classes, mas considerando casos reais ou protótipos (BOOCH et al., 2000). Sendo assim um diagrama de objetos consiste numa instância do diagrama de classes, no qual para cada classe existe um objeto (sua instância) em um determinado ponto do tempo (MELO, 2004). A figura 2.4 mostra um exemplo de diagrama de objetos.

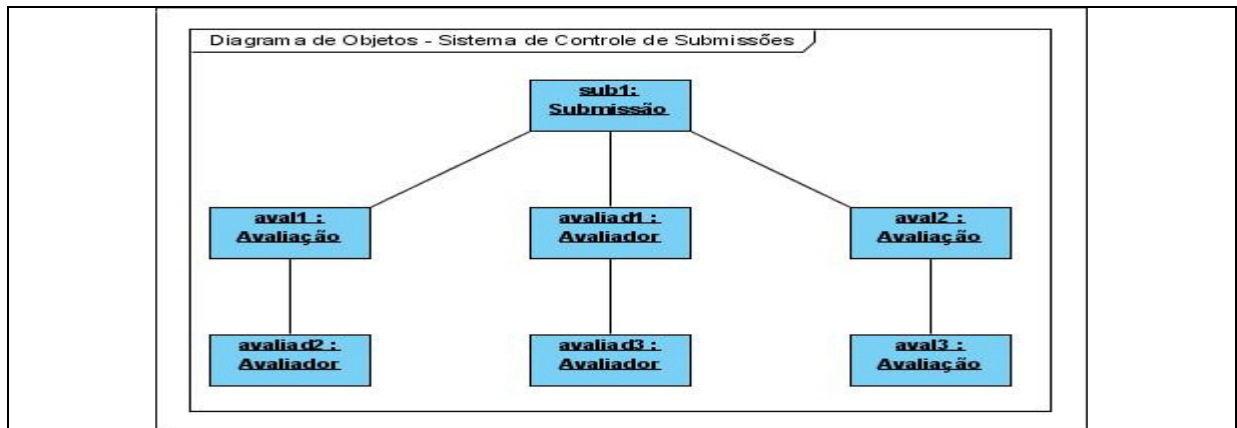


Figura 2.4 – Diagrama de Objetos (GUEDES, 2005)

### 2.2.5 Diagrama de Atividades

Um diagrama de atividades mostra o fluxo de uma atividade para outra em um sistema. Eles mostram o controle de fluxo de atividade para atividade no sistema, quais atividades estão sendo feitas em paralelo e quais caminhos alternar através do fluxo. Estes diagramas são importantes principalmente para se fazer a modelagem da função de um sistema. Os diagramas de atividades dão ênfase ao fluxo de controle entre objetos (BOOCH et al., 2000; QUATRANI, 2001). A figura 2.5 mostra um exemplo de diagrama de atividades.

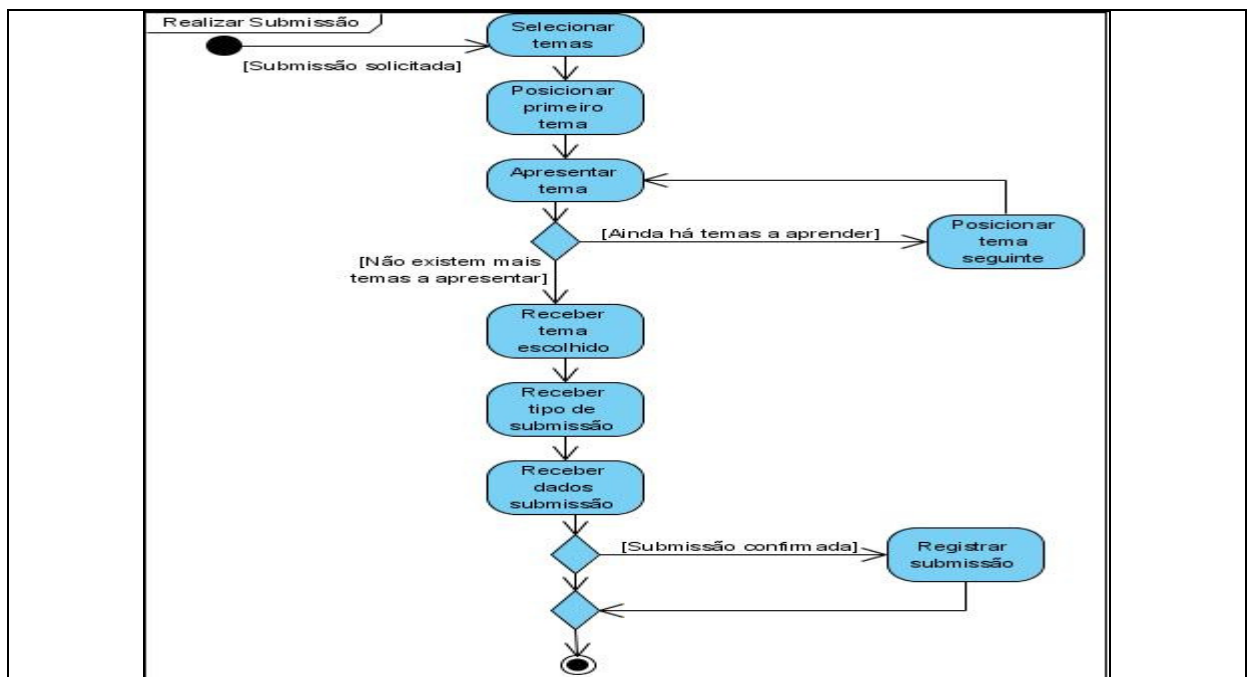


Figura 2.5 – Diagrama de Atividades (GUEDES, 2005)

## 2.2.6 Diagrama de Componentes

O diagrama de componentes identifica os componentes que fazem parte de um sistema, um subsistema ou mesmo os componentes ou classes internas de um componente individual. Um componente é identificado como sendo uma unidade modular com *interfaces* bem definidas e substituíveis dentro de seu ambiente, sendo assim pode representar tanto um componente lógico (um componente de negócio ou de processo) ou um componente físico, como arquivos contendo código fonte, arquivos de ajuda, bibliotecas, arquivos executáveis, etc. Um diagrama de componentes pode ser utilizado como uma forma de documentar como estão sendo estruturados os arquivos físicos de um sistema, permitindo uma melhor compreensão do mesmo, além de facilitar a reutilização de código (GUEDES, 2005; MELO, 2004).

Na versão 2.0 da *UML* a abordagem do diagrama de componentes foi alterada, a versão anterior tratava o diagrama de componentes como sendo um diagrama de classificadores e os artefatos que eles implementavam, este conceito está relacionado agora a definição de artefatos. A notação gráfica do diagrama também foi alterada, sendo que na versão anterior um componente era representado por um retângulo com dois retângulos menores sobressaindo-se á sua esquerda. Na versão 2.0 este símbolo foi substituído por um retângulo contendo internamente o antigo símbolo (GUEDES, 2005; MELO, 2004). A figura 2.6 representa um diagrama de componentes.

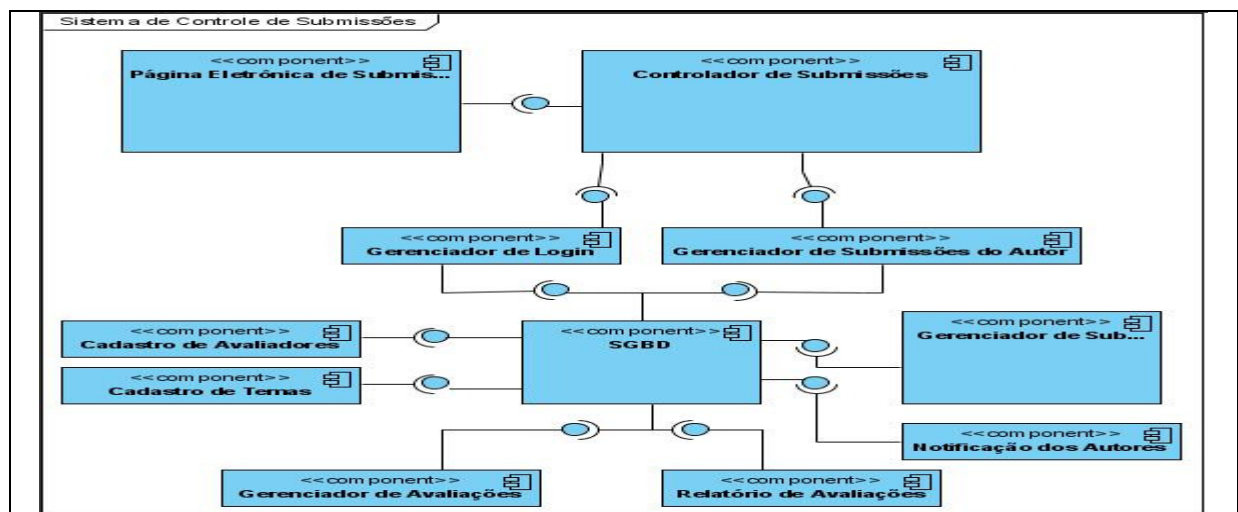


Figura 2.6 – Diagrama de Componentes (GUEDES, 2005)

## 2.2.7 Diagrama de Implantação

Um diagrama de implantação mostra a configuração dos nós de processamento em tempo de execução e os componentes neles existentes. Os elementos de um diagrama de implantação são os nós e as conexões. Um nó é uma unidade física que representa um recurso computacional e normalmente possui uma memória e alguma capacidade de armazenamento. Os nós são ligados uns aos outros através de conexões. As conexões mostram mecanismos de comunicação entre os nós: meios físicos ou protocolos de comunicação (BOOCH et al., 2000; BEZERRA, 2002; FOWLER; SCOTT, 2000). A figura 2.7 mostra um exemplo de diagrama de implantação.

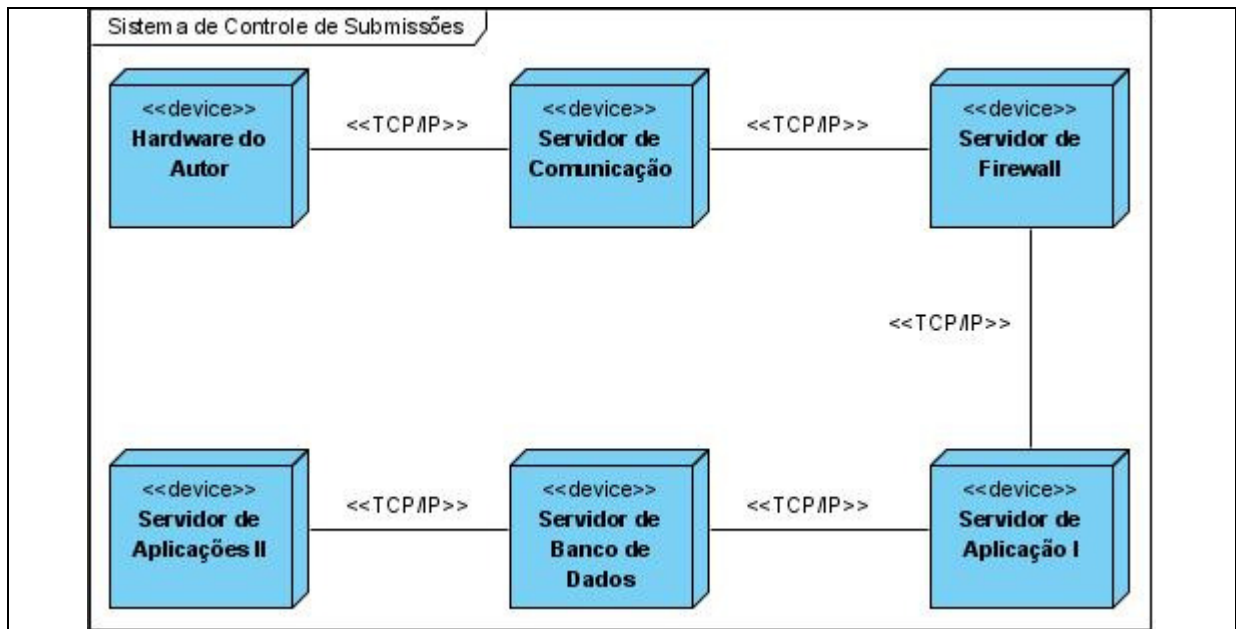


Figura 2.7 – Diagrama de Implantação (GUEDES, 2005)

### 2.2.8 Diagrama de Comunicação

O diagrama de comunicação era conhecido como diagrama de colaboração até a versão 1.5 da *UML*, tendo seu nome modificado para diagrama de comunicação a partir da versão 2.0 (GUEDES, 2005).

Um diagrama de comunicação é um diagrama de interação que dá ênfase à organização estrutural dos objetos que enviam e recebem mensagens (BOOCH et al., 2000). A colaboração entre objetos é representada por uma ligação simples acompanhada de uma numeração seqüencial e de outras informações como condições e iterações (MELO, 2004). A figura 2.8 mostra um exemplo de diagrama de comunicação.

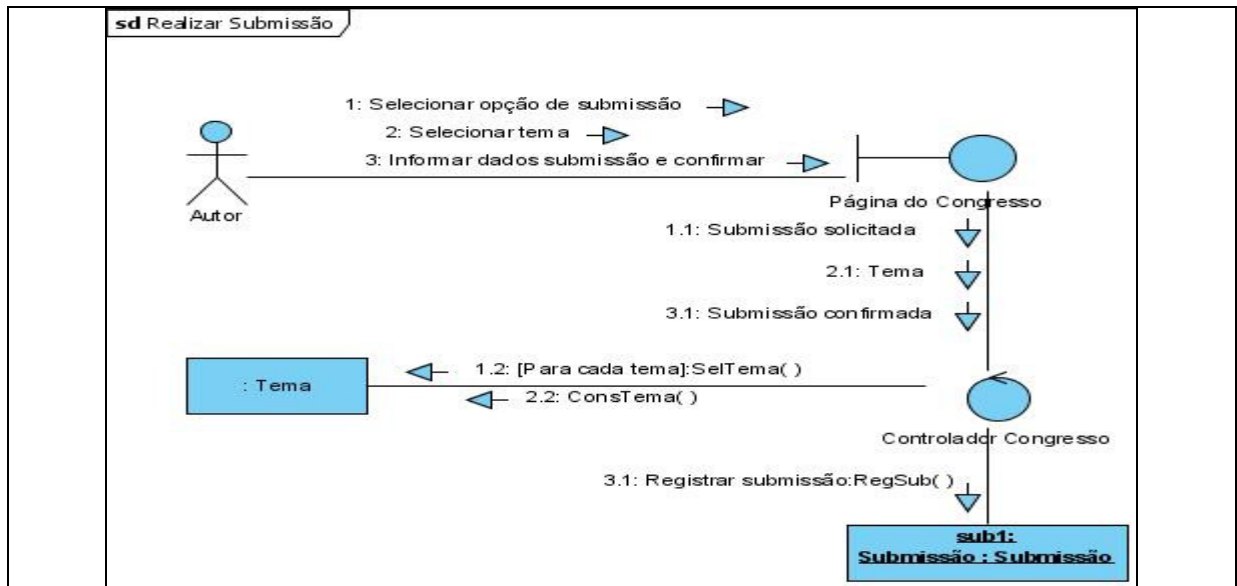


Figura 2.8 – Diagrama de Comunicação (GUEDES, 2005)

### 2.2.9 Diagrama de Máquina de Estados

Este diagrama demonstra o comportamento de um elemento através de um conjunto de transições de estado (GUEDES, 2005). Este diagrama é usado para ilustrar a visão dinâmica de um sistema. Os diagramas de máquina de estados são importantes principalmente para se fazer a modelagem do comportamento de uma *interface*, classe ou colaboração. Eles dão ênfase ao comportamento de um objeto, solicitado por eventos, que é de grande ajuda para a modelagem de sistemas reativos (BOOCH et al., 2000). A figura 2.9 traz um exemplo deste diagrama.

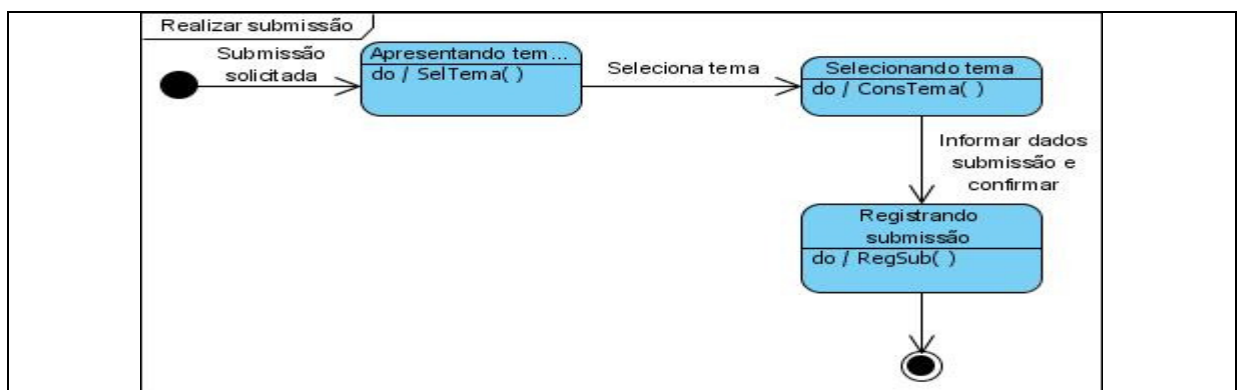
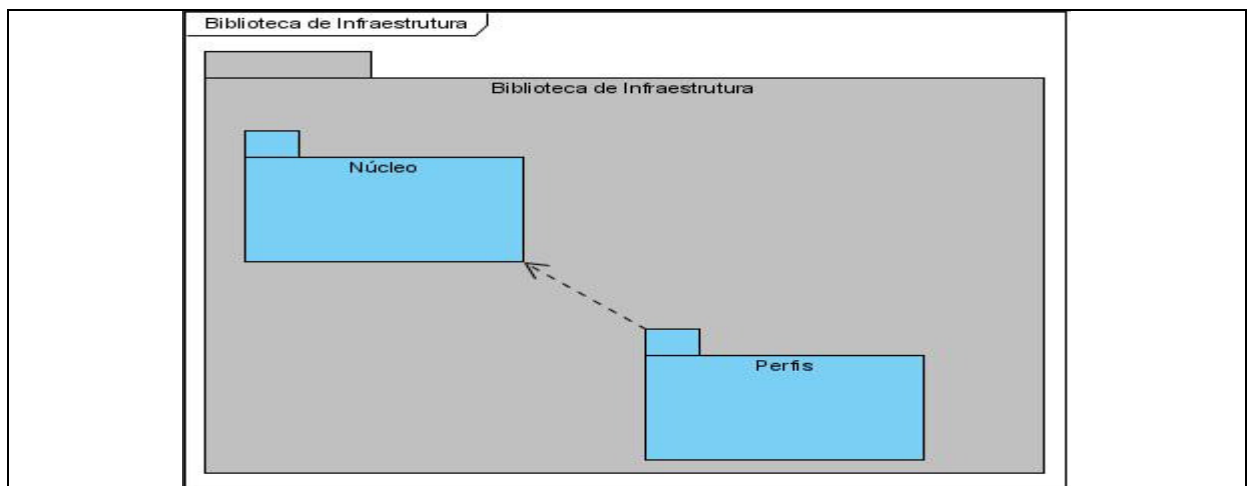


Figura 2.9 – Gráfico de Máquina de Estados (GUEDES, 2005)

### 2.2.10 Diagrama de Pacotes

A *UML* 2.0 oficializou este diagrama. Nas versões anteriores, quando era necessário mostrar pacotes e suas dependências era utilizado um diagrama de classes para fazer isto (MELO, 2004).

Um pacote é um mecanismo de propósito geral para a organização de elementos em grupos. Um pacote é representado graficamente como uma pasta com uma guia (BOOCH et al., 2000; BEZERRA, 2002). A figura 2.10 mostra um exemplo de Diagrama de Pacotes.



**Figura 2.10 – Diagrama de Pacotes (GUEDES, 2005)**

### 2.2.11 Diagrama de Estrutura Composta

Este é mais um dos diagramas incluídos na versão 2.0 da *UML*. O objetivo deste diagrama, segundo Melo (2004), é permitir que se exiba um pequeno diagrama de classes dentro de uma classe (também pode ser usado com outros classificadores). Este tipo de notação leva a uma apresentação menos confusa de um relacionamento de composição. Guedes (2005) amplia o conceito deste diagrama, segundo ele o diagrama de estrutura composta é utilizado para modelar colaborações. O termo estrutura deste diagrama refere-se a uma composição de elementos interconectados, representando instâncias de tempo de execução colaborando através de vínculos de comunicação para atingir algum objetivo comum. A figura 2.11 mostra um exemplo deste diagrama.

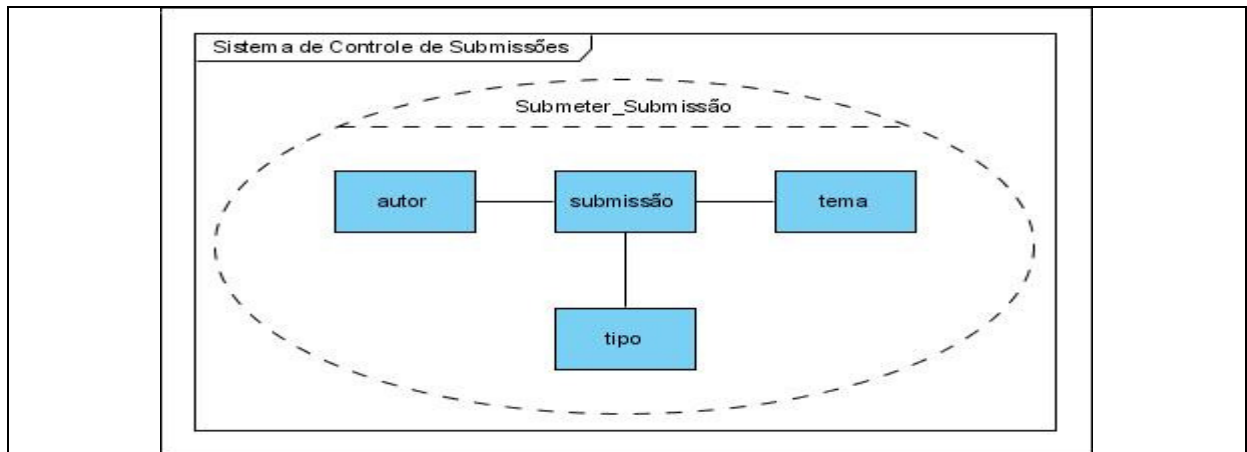


Figura 2.11 – Diagrama de Estrutura Composta (GUEDES, 2005)

### 2.2.12 Diagrama de Visão Geral

Este diagrama incluído na versão 2.0 consiste numa especialização do diagrama de atividades representando interações, promovendo uma visão geral do fluxo de controle. Este diagrama utiliza as notações gráficas do diagrama de atividades, onde os nós (atividades) ou são interações ou são ocorrências de interações (MELO, 2004; GUEDES, 2005). O exemplo a seguir, figura 2.12, mostra um diagrama de visão geral.

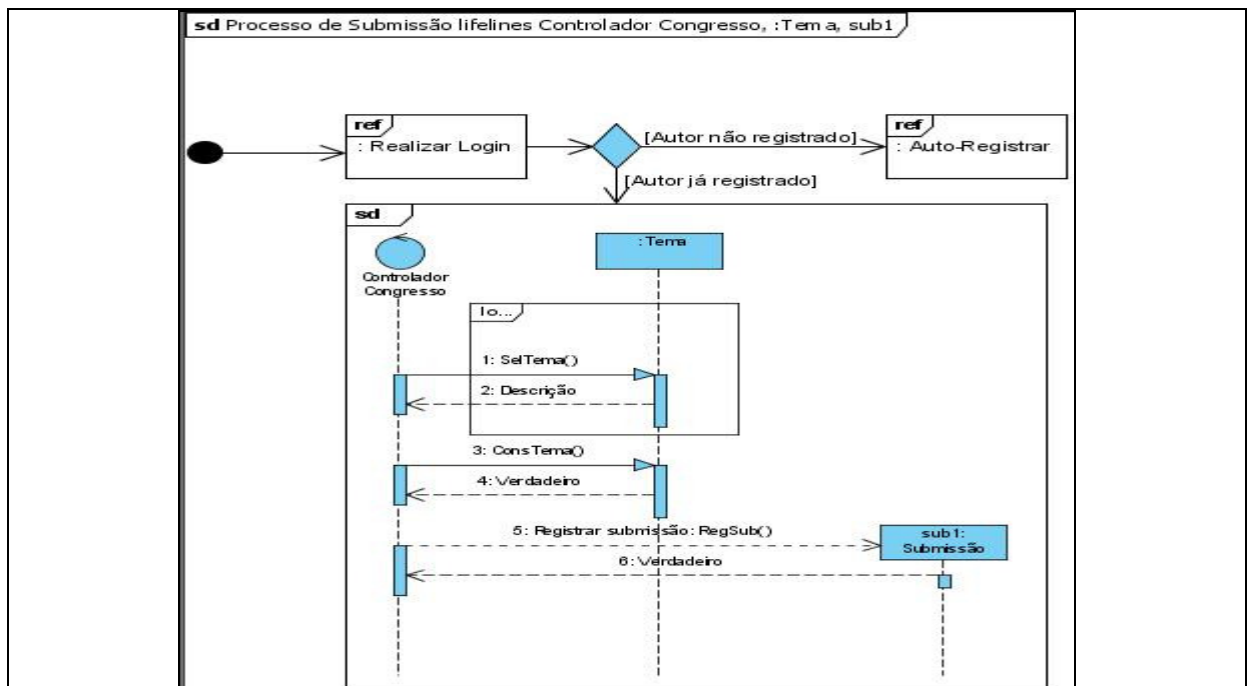
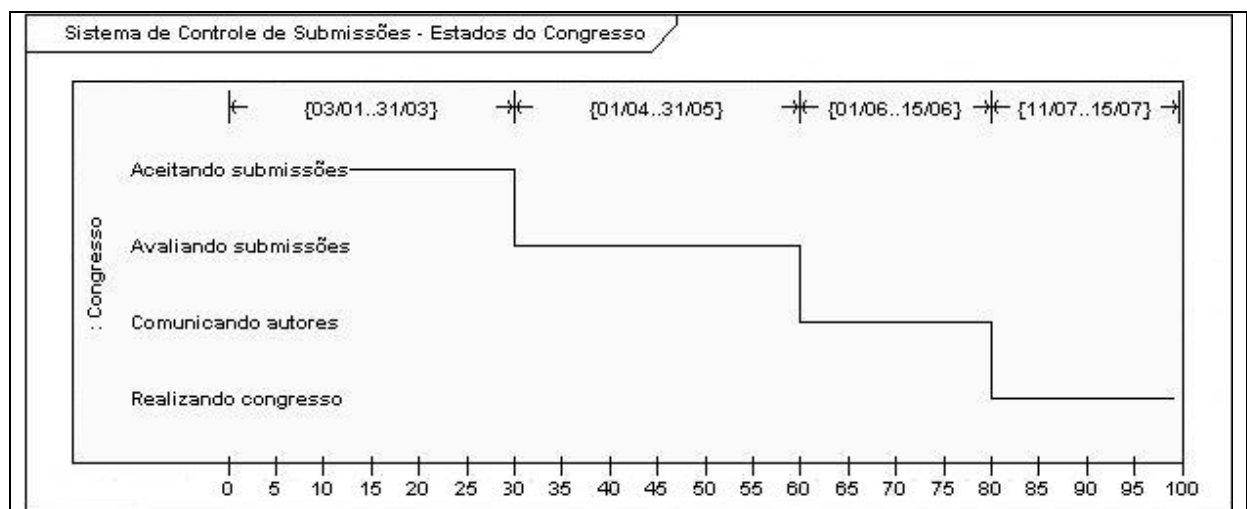


Figura 2.12 – Diagrama de Visão Geral (GUEDES, 2005)

### 2.2.13 Diagrama Temporal

O diagrama de tempo ou temporal mostra as mudanças de um objeto ao longo do tempo, em resposta a eventos externos. Este diagrama terá pouca utilidade para modelar aplicações comerciais, porem poderá ser utilizado na modelagem de sistemas de tempo real ou sistemas que utilizem recursos multimídia/hipermídia, onde o tempo em que um objeto executa algo é muitas vezes importante (GUEDES, 2005; MELO, 2004). A figura 2.14 mostra um exemplo deste diagrama.



**Figura 2.13 – Diagrama Temporal (GUEDES, 2005)**

No capítulo seguinte são abordados alguns conceitos sobre a utilização de *frameworks* orientados a objetos.



### **3 FRAMEWORKS ORIENTADOS A OBJETOS**

Este capítulo aborda os conceitos a respeito da utilização de *frameworks*. Os *frameworks* estão cada vez mais presentes na área de desenvolvimento de *software*, pois são a maneira pela qual os sistemas orientados a objetos conseguem a maior reutilização (GAMMA et. al, 2000).

Na literatura podem ser encontradas muitas definições para *frameworks* orientados a objetos, a seguir são apresentadas algumas delas.

Segundo Silva (2000) à abordagem de *frameworks* orientados a objetos utiliza o paradigma de orientação a objetos para produzir uma descrição de um domínio para ser reutilizada. Sendo assim um *framework* é uma estrutura de classes inter-relacionadas, que correspondem a uma implementação incompleta para um conjunto de aplicações de um domínio, sendo que esta estrutura de classes deve ser adaptada para a geração de aplicações específicas.

Para Pree (1995) um *framework* é uma coleção de classes abstratas e concretas que representam um subsistema, estas classes (abstratas e concretas) podem ser estendidas ou adaptadas para construir um novo subsistema.

Para D'Souza e Wills (1998) um *framework* pode ser interpretado como sendo um pacote de *templates*, um pacote desenvolvido para ser importado com substituições. Ele é entendido para prover uma nova versão, baseada nas substituições realizadas.

Segundo Gamma et al. (2000) um *framework* predefine os parâmetros de projeto, de maneira que o projetista/implementador da aplicação possa se concentrar nos aspectos específicos da sua aplicação.

Segundo Mattsson (1996) *framework* é uma arquitetura desenvolvida com o objetivo de se obter a máxima reutilização, representada como um conjunto de classes abstratas e concretas, com grande potencial de especialização.

### **3.1 Framework X Bibliotecas**

É importante salientar a diferença entre *frameworks* e bibliotecas de classes. *Frameworks* diferem de bibliotecas porque apresentam uma inversão do fluxo de controle, entre os *frameworks* e os desenvolvedores dos *frameworks*. Projetos baseados em bibliotecas forçam o desenvolvedor a reescrever todo o fluxo de controle da aplicação, enquanto projetos baseados em *frameworks* colocam o *framework* no controle do fluxo, onde o desenvolvedor reescreve apenas métodos especiais que adaptam alguma funcionalidade do *framework*, os quais são chamados pelo próprio *framework* (GAMMA, 2000; JOHNSON; FOOTE, 1988).

### **3.2 Framework Model View Controller (MVC)**

Um dos primeiros *frameworks* que se tem conhecimento é o *MVC* desenvolvido no ambiente do *Smalltalk-80* (PREE, 1995). Este *framework* é muito utilizado hoje em dia na comunidade de *software*.

Através do *MVC* é possível dividir uma aplicação em três camadas distintas, desta forma separando a lógica da aplicação chamada de *Model* (modelo), da *interface* do usuário conhecida como *View* (visão) e do fluxo da aplicação, o *Controller* (controlador).

Cada camada anteriormente citada corresponde a uma classe abstrata do *framework MVC*, e elas cooperam através de um protocolo de interação bem definido. Essas classes são definidas da seguinte maneira (GERBER, 1999):

- *Model*: um objeto do tipo *Model* pode manter objetos dependentes (*views* e *controllers*), e enviar mensagens para notificar modificações no modelo para os dependentes. A classe *Model* define um protocolo de mensagens específico para estes propósitos. Dessa forma um modelo de uma aplicação é considerado como sendo uma instância de uma subclasse de *Model*, que deve implementar o comportamento específico da aplicação;

- *View*: um objeto *View* pode ser dividido em outros objetos do tipo *View*, chamados de subvisões. A classe *View* define um protocolo para interagir com objetos do tipo *controller* e *model* e interagir com suas subvisões, desta forma coordenando ações de transformação e apresentação. Um objeto *View* em uma aplicação é uma instância de uma subclasse de *View*. A subclasse tem que fornecer a implementação para o protocolo de apresentação, para definir detalhes específicos do modelo que representa;
- *Controller*: esta classe define os protocolos para manipular os objetos do tipo *Model* e *View* a partir de mensagens de interação enviadas pelo usuário através de determinados dispositivos de entrada de dados.

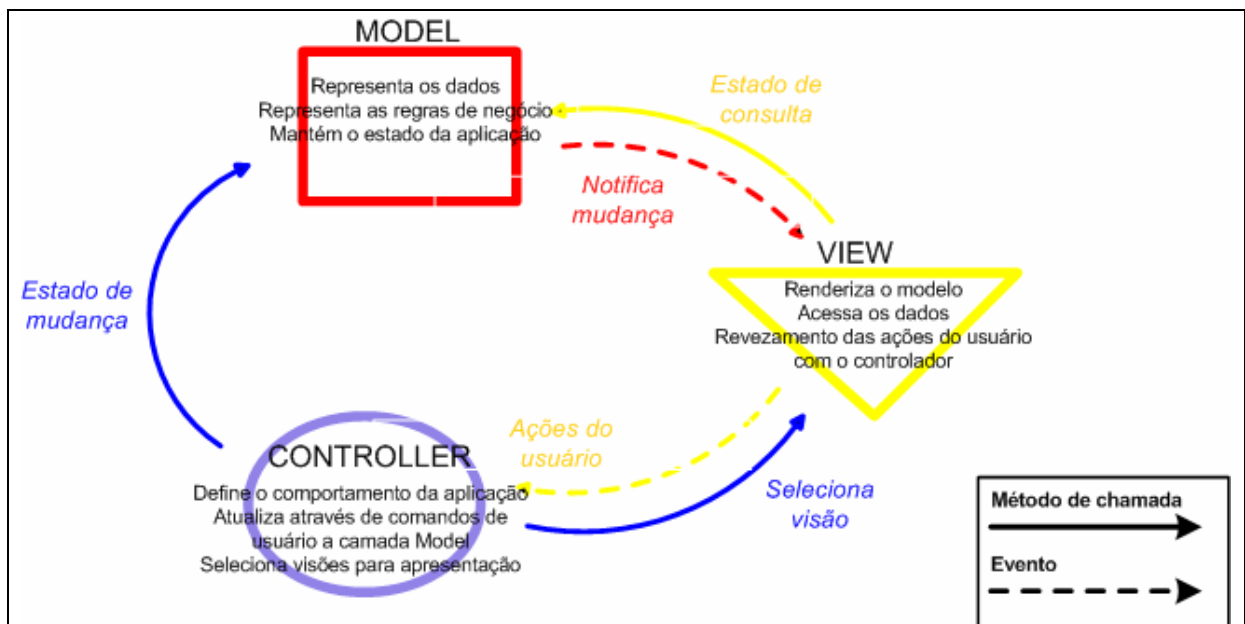


Figura 3.1 - MVC

### 3.3 Classificação dos Frameworks

A classificação de *frameworks* se dá de várias formas e de vários pontos de vista, segundo alguns autores, caracterizando tipos distintos de *frameworks* de acordo com seus propósitos e de acordo com as técnicas usadas para entendê-los.

### 3.3.1 Quanto ao Propósito

Fayad e Schmidt dividem este tipo de *framework* em infra-estrutura, integração e domínio de aplicação (FAYAD; SCHMIDT, 1997).

*Frameworks* de Infra-estrutura: são usados para dar suporte a diversos tipos de aplicações, independentemente dos domínios endereçados por estas aplicações, desta forma simplificando o desenvolvimento da infra-estrutura. Exemplos conhecidos deste *framework* seriam *frameworks* para sistemas operacionais, sistemas de comunicação, *interfaces* com o usuário e ferramentas para processamento de linguagens naturais;

*Frameworks* de Integração: são usados normalmente para integrar aplicações e componentes distribuídos. *Frameworks* de integração são desenvolvidos também para aumentar as habilidades de desenvolvedores de *software* em modularizar, reusar e estender a infra-estrutura do *software* para trabalhar em um ambiente distribuído. Exemplos comuns deste *framework* incluem o *ORB (Object Request Broker) Frameworks* e bases de dados transacionais;

*Frameworks* de Domínio de Aplicação: representam esqueletos de aplicações para domínios específicos, como por exemplo, para telecomunicações, manufatura e engenharia financeira. O desenvolvimento deste tipo de *framework* geralmente é caro, assim como sua compra. Entretanto eles podem prover um retorno substancial do investimento, quando utilizados para o suporte ao desenvolvimento de aplicações para usuário final e para a criação de produtos em domínios específicos. *Frameworks* deste tipo capturam elementos invariantes de domínio, e deixam em aberto os aspectos específicos de cada aplicação.

Algumas diferenças básicas entre *Frameworks* de Domínio de Aplicação e *Frameworks* de Infra-estrutura e Integração está no fato que estes dois preocupam-se basicamente com problemas internos de desenvolvimento de *software* e são independentes do domínio de aplicação. Entretanto esses *frameworks* são essenciais para o desenvolvimento de *softwares* de qualidade, porém não geram um retorno significativo do investimento. Para as empresas, este tipo de *framework* é mais vantajoso ser adquirido do que desenvolvido internamente.

### 3.3.2 Quanto as Técnicas de Reuso

De acordo com esse escopo os *frameworks* são divididos em *Frameworks* de Caixa Branca (*whitebox*), *Frameworks* de Caixa Preta (*blackbox*) e *Frameworks* de Caixa Cinza (*graybox*)(FAYAD; SCHMIDT, 1997).

Em *Frameworks* de Caixa Branca o desenvolvedor deve estender as classes que constituem o *framework* para criar aplicações específicas, isto é, a reutilização é obtida através da herança. Isso dificulta muito a utilização, pois o desenvolvedor precisa entender como funciona cada detalhe do *framework* para poder utilizá-lo. Dessa forma, ensinar como funciona este tipo de *framework* é quase a mesma coisa que ensinar a criá-lo (JOHNSON; FOOTE, 1988; FAYAD; SCHMIDT, 1997).

Em um *Framework* de Caixa Preta o desenvolvedor combina diversas classes concretas existentes do *framework* para obter uma aplicação, desta forma o reuso é obtido através das composições feitas com estas classes. Neste caso o desenvolvedor precisa apenas entender as *interfaces* para poder utilizá-lo (JOHNSON; FOOTE, 1988; FAYAD; SCHMIDT, 1997).

Johnson afirma que o ideal é que cada *framework* de caixa branca evolua para um *framework* de caixa preta com o tempo. Dessa forma ele diz que *frameworks* de caixa branca devem ser vistos como sendo uma fase natural na evolução de um sistema, que será algum dia um passo na evolução de uma coleção de métodos em um conjunto de componentes (JOHNSON; FOOTE, 1988).

Um *Framework* de Caixa Cinza é considerado como sendo uma mistura dos conceitos dos *frameworks* de Caixa Branca e Caixa Preta. Dessa forma o reuso é obtido por meio da herança, por ligação dinâmica e também pelas *interfaces* de definição (BRAGA, 2002).

Um aspecto variável de um domínio de aplicação é chamado de ponto variável (*hotspot*) (BUSCHMAN et al., 1996). Diferentes aplicações dentro de um mesmo domínio são diferenciadas por um ou mais pontos variáveis. Eles representam as partes do *framework* de aplicação que são específicas de cada sistema. Os pontos variáveis são projetados para serem genéricos e podem ser adaptados às necessidades da aplicação (BRAGA, 2002).

Pontos fixos (*frozen-spots*) definem a arquitetura geral de um sistema de *software*, seus componentes básicos e os relacionamentos entre eles. Estes pontos permanecem sem nenhuma modificação em todas as instanciações de um *framework* de aplicação (BRAGA, 2002).

Segundo um estudo intitulado “*Application Frameworks: A Survey*” realizado por Yassin e Fayad (FAYAD; JOHNSON, 1999), a tendência de desenvolvimento é pelos *frameworks* de caixa cinza. Apesar dos *frameworks* de caixa branca serem mais fáceis de desenvolver, a exposição do código fonte aos desenvolvedores pode causar problemas. Já os de caixa preta são mais abstratos e menos flexíveis, dificultando desta forma a manutenção de aplicações derivadas dele. *Frameworks* de caixa cinza quando bem projetados superam as barreiras impostas pelos outros dois tipos, porque possuem alta flexibilidade e uma grande facilidade de extensão, sem expor desnecessariamente informações que não interessam aos desenvolvedores (BRAGA, 2002).

Em resumo, um *framework* caixa branca é muito mais fácil de projetar, pois não existe a necessidade de se prever todas as alternativas de implementação possíveis ao contrário dos de caixa preta, que necessitam dessa previsão. No entanto os *frameworks* de caixa preta são mais fáceis de utilizar, pois basta apenas escolher a implementação desejada, enquanto os de caixa branca é necessário fornecer um implementação completa. *Frameworks* de caixa cinza tentam tirar o máximo proveito das vantagens dos outros dois tipos de *framework*, sendo assim, considerado como um meio termo. *Frameworks* de caixa branca podem evoluir para se tornar cada vez mais caixa preta (JOHNSON, 1997). Isto pode ser alcançado de forma gradativa, implementando-se várias alternativas que posteriormente podem ser utilizadas na instanciação do *framework* (BRAGA, 2002).

### **3.4 Compreensão e Uso de *Frameworks***

Utilizar um *framework* consiste em criar aplicações a partir deste *framework*. A motivação para se usar um *framework* na construção de aplicações é a perspectiva de um aumento considerável da produtividade e da qualidade do sistema desenvolvido, em função da reutilização promovida. Para produzir uma aplicação utilizando um *framework* como base, deve-se estender sua estrutura de modo a cumprir os requisitos básicos da aplicação a que se quer desenvolver (SILVA, 2000).

Um dos obstáculos encontrados para se desenvolver aplicações através do uso de *frameworks* é justamente compreender seu uso, como funciona e interage internamente. A questão é que se um desenvolvedor de aplicações necessitar empreender um esforço considerável para aprender a utilizar um *framework*, que para desenvolver uma aplicação do início, a utilidade deste *framework* acaba se tornando questionável (SILVA, 2000).

Dessa forma torna-se imperativo do processo de desenvolvimento de *frameworks* a produção de subsídios para minimizar o esforço do desenvolvedor em compreender como utilizar os *frameworks* produzidos. Johnson classifica diferentes níveis de compreensão de um *framework*: saber que tipos de *softwares* podem ser desenvolvidos através da utilização do *framework* (a qual domínio de aplicações ele pertence), ser capaz de desenvolver aplicações elementares sob o *framework* e conhecer os detalhes a respeito do projeto do *framework*. Para capacitar ao desenvolvimento de aplicações elementares recomenda o desenvolvimento de uma documentação no estilo *cookbook*<sup>8</sup>, que não entre em detalhes de projeto (JOHNSON, 1994). A utilização de *frameworks* a partir de *cookbooks* torna-se limitada quanto a gama de aplicações que podem ser desenvolvidas: se um determinado requisito de um aplicação não é tratado no conjunto de instruções, o mecanismo não consegue auxiliar o desenvolvimento. Para suprir tal deficiência, usuários de *frameworks*, invariavelmente, são forçados a buscar conhecimento sobre o projeto do *framework*. Dessa forma, é extremamente importante que a descrição de um *framework* fornecida a um usuário o habilite a desenvolver aplicações com o menor esforço possível, mas que também forneça os meios de se aprofundar a compreensão do projeto do *framework*, a fim de habilitá-lo a produzir aplicações mais complexas, sem que lhe seja exigido um esforço excessivo. O fator de comparação diz que o esforço compreendido para desenvolver uma aplicação através de um *framework* não deve ser maior que o necessário para produzir a mesma aplicação sem utilizar o *framework* (SILVA, 2000; JOHNSON, 1994).

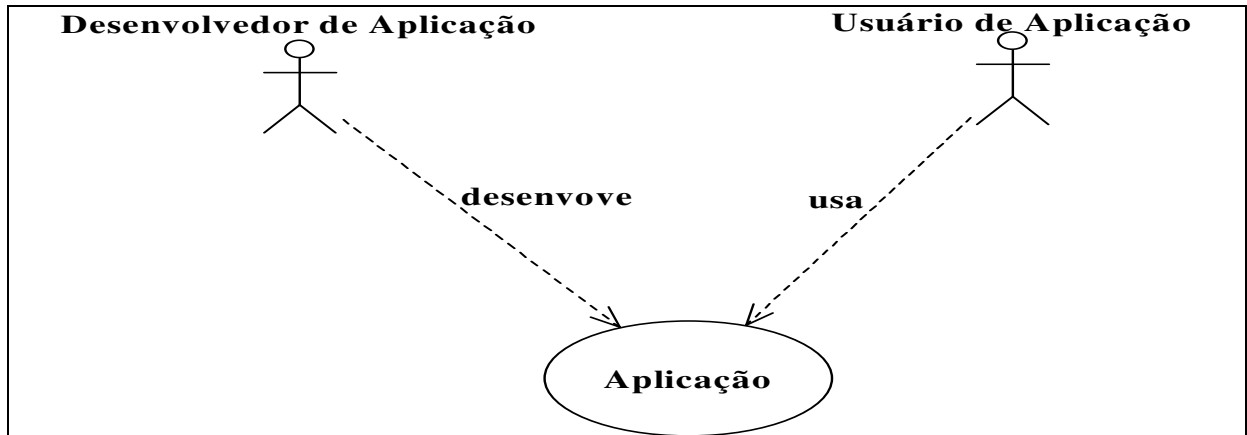
### 3.5 Componentes Envolvidos no Desenvolvimento e Uso de *Frameworks*

O desenvolvimento tradicional de aplicações envolve dois tipos de personagens: o desenvolvedor de aplicação e o usuário de aplicação. A tarefa dos desenvolvedores consiste em levantar os requisitos de uma determinada aplicação, desenvolvê-la (o que inclui toda a

---

<sup>8</sup> Livro de Receitas - Documentação que mostra como utilizar um *framework* através de exemplos passo a passo.

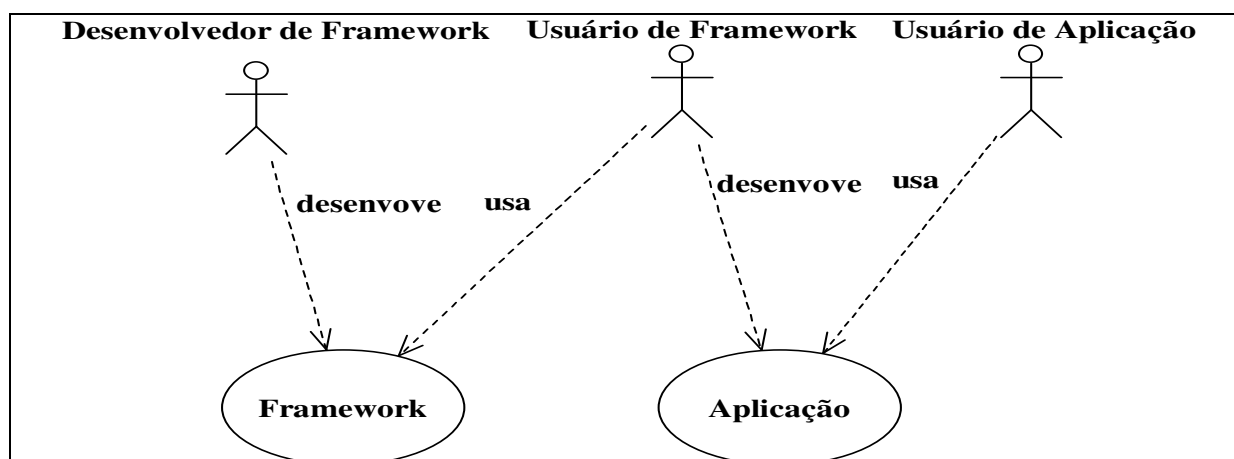
documentação que ensina a utilizar a aplicação) e entregá-la aos usuários. Os usuários por sua vez interagem com a aplicação apenas através da *interface* que lhes é oferecida. A figura 3.2 exemplifica essa relação entre desenvolvedores e usuários de uma aplicação (SILVA, 2000).



**Figura 3.2 – Componentes do Desenvolvimento tradicional de aplicações (SILVA, 2000)**

O desenvolvimento de *frameworks* introduz outro componente, além do desenvolvedor e usuário de aplicação: o desenvolvedor de *framework*. No contexto dos *frameworks*, o papel do usuário de aplicação permanece o mesmo. O papel do desenvolvedor de aplicações difere do caso anterior pela inserção do *framework* no processo de desenvolvimento de aplicações. Dessa forma o desenvolvedor de aplicações se torna um usuário do *framework*, que deve estender e adaptar a estrutura deste *framework* para a produção de aplicações. Ele permanece com as mesmas funções do caso anterior: obter os requisitos da aplicação, desenvolve-la usando o *framework* (o que em geral, não dispensa o desenvolvedor de aplicações da necessidade de produzir código) e entregá-la aos usuários. O novo componente criado no contexto dos *frameworks*, o desenvolvedor de *framework*, tem a responsabilidade de produzir *frameworks* e algum modo de ensinar os usuários do *framework* a desenvolver aplicações. A figura 3.3 exemplifica essa nova relação (SILVA, 2000).





**Figura 3.3 – Componentes do desenvolvimento de aplicações baseadas em *frameworks* (SILVA, 2000)**

### 3.6 Construção de *Frameworks*

Na construção de um *framework* diversas atividades não sequenciais se repetem, até se obter uma estrutura de classes que contemple os requisitos de generalidade, flexibilidade e extensibilidade. Para uma melhor compreensão podem-se organizar estas atividades que constituem o processo de desenvolvimento da estrutura de classes de um *framework* nas seguintes etapas: etapa de generalização, etapa de flexibilização, etapa de aplicação de metapadrões, etapa de aplicações de padrões de projeto e etapa de aplicação de princípios práticos de orientação a objetos (SILVA, 2000).

A etapa de generalização consiste no processo de identificar estruturas idênticas nas aplicações analisadas, e na fatoração de estruturas em que identificam semelhanças, de modo a obter uma estrutura de classes que generalize o domínio tratado. Para a busca de generalidade da estrutura do *framework*, um dos primeiros aspectos é a identificação das classes que representam elementos e conceitos do domínio tratado. Isto pode ser obtido através da confrontação de estruturas de classes de diferentes aplicações. Nesta etapa também são pesquisados elementos prontos que podem ser reutilizados, como por exemplo, outros *frameworks* e componentes (SILVA, 2000; BRAGA, 2002).

Na etapa de flexibilização o principal objetivo é identificar o que deve ser mantido flexível na estrutura de classes que generaliza o domínio, de modo que a estrutura possa ser especializada, gerando diferentes aplicações. É um processo de localização de pontos variáveis na estrutura de classes e de escolha de soluções de projeto para modelá-los. A

identificação destes pontos variáveis ocorre quando se determinam situações de processamento comuns às aplicações do domínio (os pontos fixos), porém, com variações de uma aplicação para outra (SILVA, 2000; BRAGA, 2002).

Quando identificado um ponto variável, deve-se comparar o requisito de flexibilização por ele imposto com os casos tratados por padrões existentes. Se um padrão é considerado adequado para a solução do problema, ele pode ser incorporado a estrutura de classes do *framework*. No caso específico dos meta-padrões, o seu uso consiste em transformar um procedimento geral em um método *template*, cujo comportamento é flexibilizado através da dependência de métodos *hook* (gancho), que podem ter diferentes implementações (SILVA, 2000).

A etapa de aplicação de padrões de projeto consiste em incluir as classes de um padrão selecionado na estrutura de classes do *framework* (ou fazer com que classes preexistentes assumam responsabilidades correspondentes as classes do padrão de projeto)(SILVA, 2000; BRAGA, 2002).

Segundo Johnson e Foote (1988) para se desenvolver um *framework* com uma estrutura de classes flexíveis são necessários seguir alguns princípios básicos da orientação a objetos, como por exemplo, o uso da herança para a reutilização de *interfaces*, ao invés do uso de herança para a reutilização de código. A reutilização de código através de composição de objetos, a preocupação em promover polimorfismo na definição de classes e métodos, de modo a possibilitar acoplamento dinâmico.

### **3.7 Documentação de *Frameworks***

*Frameworks* são projetos reusáveis e não apenas linhas de código, são mais complexos de se desenvolver do que simples aplicações, tornando sua documentação mais difícil de ser realizada.

A documentação de *frameworks* deve atender a diversos requisitos. Estes requisitos podem ser alcançados fazendo uma estruturação da documentação através do uso de padrões ou de alguma linguagem de padrões. O principal propósito deste conjunto de padrões é mostrar como utilizar um *framework*, não apenas mostrar como ele funciona. Com a utilização de padrões pode-se descrever o propósito de um *framework*, permitindo aos

usuários do *framework* desenvolverem sem ter a necessidade de conhecer os detalhes internos de funcionamento (RÉ, 2002).

Geralmente as documentações dos *frameworks* descrevem primeiramente como um *framework* trabalha e em seguida como utilizá-lo. No entanto, não é possível entender completamente a estrutura de um *framework* até que ele tenha sido utilizado efetivamente. A utilização de exemplos é uma boa forma de comprovar se a forma de utilização foi ou não bem entendida pelo usuário, além de tornar os *frameworks* mais completos e facilitarem o entendimento do fluxo de controle (JOHNSON, 1992).

### **3.8 Vantagens e Desvantagens do Uso de *Frameworks***

Segundo Fayad e Schmidt (1997) os benefícios primários da utilização de *frameworks* são os seguintes:

**Modulariedade:** *frameworks* aumentam a modulariedade através de um encapsulamento de detalhes de implementação flexíveis, sob *interfaces* estáveis. A modulariedade do *framework* auxilia no aumento da qualidade do *software* localizando lugares de impacto no projeto e os impactos que mudanças na implementação podem causar. Esta localização reduz o esforço necessário para compreender e manter um *software* existente.

**Reusabilidade:** as *interfaces* estáveis providas pelo *framework* aumentam a reusabilidade através da definição de componentes genéricos que podem ser replicados para criar novas aplicações. A reusabilidade fornecida pelo *framework* permite que desenvolvedores utilizem seus conhecimentos no domínio do problema para evitar a recriação e a revalidação de soluções já conhecidas, presentes no *framework*. O reuso de componentes da estrutura do *framework* pode render melhorias substanciais na produtividade dos desenvolvedores, aumentando a qualidade, o desempenho, a confiabilidade e a interoperabilidade do *software*.

**Extensibilidade:** Provê métodos explícitos que possibilitam às aplicações estenderem suas *interfaces*. A extensibilidade do *framework* é essencial para assegurar uma customização oportuna de serviços e características novas em uma aplicação.

Inversão de Controle: a arquitetura de um *framework* é caracterizada pela inversão de controle. Esta arquitetura permite que a aplicação possa processar eventos que podem ser customizados e chamados pelo próprio *framework*. Quando estes eventos ocorrem o *framework* reage invocando métodos pré-definidos de cada aplicação, que executam processos pré-definidos nos eventos.

Segundo Mattsson (1996) as principais desvantagens no uso e construção de *frameworks* envolvem:

A dificuldade de desenvolver bons *frameworks*, pois é necessário ter um alto conhecimento no domínio do *framework* para poder construí-lo de forma que atenda as necessidades.

A documentação do *framework* é crucial para o desenvolvedor de aplicações. Se esta documentação não estiver bem elaborada irá dificultar o uso do *framework* para desenvolvimento de novas aplicações.

O processo de depuração de aplicações originadas através de um *framework* pode ser muito complicado, pois é muito difícil distinguir se um erro está no código do *framework* ou no da aplicação. Se os erros estão localizados no *framework* pode ser impossível para o usuário consertá-los.

Em alguns casos a generalidade e a flexibilidade de um *framework* podem trabalhar contra ele, reduzindo a eficiência em uma aplicação em particular.

Quanto mais customizado um componente é, mais ele se adapta a uma situação em particular, mas mais esforço é necessário para aprender a usá-lo e para usá-lo efetivamente.

## 4 API'S DE MECANISMOS DE BUSCA

Este trabalho tem por objetivo primeiramente a construção de um *framework* para PLN utilizando *API's* de mecanismos de busca, e posteriormente o desenvolvimento de uma aplicação simples baseada nesse mesmo *framework*. O processamento da linguagem natural é feito sobre um conjunto de textos chamado *corpus*. Para montar este conjunto é necessário à reunião de diversos textos, para isto este trabalho propõe que este *corpus* seja montado utilizando a base de dados contida na *web*. Para se ter acesso a essa base de dados alguns mecanismos de busca disponibilizam gratuitamente suas *API's* de pesquisa na *web* para desenvolvimento de aplicações não comerciais. O presente trabalho utiliza as *API's* disponibilizadas pelo *Google*, *Yahoo* e *Technorati* para realizar pesquisas e montar o *corpus* necessário para o processamento das técnicas de PLN.

### 4.1 Google API

A *Google Web API* nada mais é do que uma licença para se utilizar a base de dados do mecanismo, desde que seja para uso pessoal, onde são disponibilizadas ferramentas e metodologias para troca de dados entre um *software* e a *API*. Esta *API* foi disponibilizada pelo mecanismo *Google* como forma de manter seu desenvolvimento e comprometimento com seus usuários, dessa forma tornando sua base disponível no futuro a este tipo de utilização, massificar seu uso e se manter como referência em inovação quando se fala em busca de informações na *web* (TAGLIASSUCHI, 2002). Para se obter essa licença ou chave de pesquisa é necessário que o desenvolvedor registre uma conta na página do *Google*.

Através desta *API* o desenvolvedor tem acesso a uma infinidade de informações contidas nas páginas que estão armazenadas na base de dados do *Google*, permitindo que aplicações possam buscar e manipular estas informações, de acordo com a vontade do desenvolvedor de *software*.

De acordo com o *Google* (2006) as suas *API's* podem ser usadas para construir diversas aplicações, deixando na mão do desenvolvedor a tarefa de criá-las de acordo com sua “genialidade”. Alguns exemplos de possibilidades citados pelo mecanismo são:

- Criar mecanismos que ficam monitorando a *web* regularmente em busca de novas informações a respeito de um determinado assunto;
- Criar mecanismos que realizam pesquisas de mercado através da análise de algumas diferenças nas informações disponibilizadas;
- Fazer pesquisas utilizando *interfaces* não *HTML*<sup>9</sup>, como linhas de comando e aplicações;
- Criar jogos inovadores utilizando as informações disponíveis na *web*.

A tabela 4.1 lista todos os parâmetros válidos que podem ser usados quando for feita uma solicitação de pesquisa através da *API* e descreve como estes parâmetros podem influenciar na pesquisa realizada.

Nome	Descrição
<i>key</i>	Chave fornecida pelo <i>Google</i> , esta chave é requerida para ter acesso ao Serviço do <i>Google</i> . O <i>Google</i> a utiliza para autenticar e logar no mecanismo.
<i>q</i>	Termos de pesquisa que serão utilizados na requisição. Ex.: “+”, “-“, “OR”, etc.
<i>start</i>	Zero baseado no índice do primeiro resultado desejado.
<i>maxResults</i>	Número de resultados desejados por pesquisa realizada. O Valor máximo é de 10 resultados.
<i>filter</i>	Ativa ou desativa o filtro de resultados. Este filtro é responsável por mostrar ou esconder informações similares no resultado.
<i>restricts</i>	Restringe a pesquisa. Ex.: Realizar pesquisas que retornem somente resultados de páginas de um determinado país ou determinado tópico.

<sup>9</sup> *HTML* – Do Inglês *Hiper Text Markup Language* (Linguagem de Formatação de Hipertexto). Linguagem de marcação utilizada para produzir páginas na Internet. De modo geral são documentos de texto escritos em códigos que podem ser interpretados pelos *browsers* para exibir as páginas da *World Wide Web*. Disponível em <http://www.w3.org/MarkUp/>.

<i>safeSearch</i>	Um valor booleano que habilita os filtros de conteúdo adulto para a pesquisa.
<i>lr</i>	<i>Language Restrict</i> – Restringe a pesquisa a documentos de uma determinada língua.
<i>ie</i>	<i>Input Encoding</i> – Funcionalidade em desuso que será retirada nas próximas versões da <i>API</i> . Todas as pesquisas devem ser feitas utilizando <i>UTF-8</i> <sup>10</sup> .
<i>oe</i>	<i>Output Encoding</i> – Funcionalidade em desuso que será retirada nas próximas versões da <i>API</i> . Todas as pesquisas devem ser feitas utilizando <i>UTF-8</i> .

**Tabela 4.1 – Parâmetros de Pesquisa do Google (GOOGLE, 2006)**

A tabela 4.2 mostra algumas limitações impostas pelo mecanismo do *Google* para a realização de pesquisas através do uso da *API*. Segundo o *Google* (2006) algumas destas limitações são decorrentes da infra-estrutura do mecanismo, que está otimizada somente para uso de usuários finais (usuários que acessam o *Google* para fazer algum tipo de pesquisa). Entretanto futuramente esses limites serão alterados.

<b>Componente</b>	<b>Limite</b>
Tamanho da requisição de pesquisa	2048 bytes
Número máximo de palavras na pesquisa	10
Número máximo do site: termos na pesquisa	1 (por cada pesquisa)
Número máximo de resultados por pesquisa	10
Valor máximo de <start> + <maxResults>	1000

**Tabela 4.2 – Limitações da API do Google (GOOGLE, 2006)**

Além das limitações quanto à realização de pesquisas, a *API* também tem algumas limitações no que diz respeito às linguagens que ela está preparada para suportar (linguagens que já foram homologadas). As linguagens que esta *API* da suporte são as seguintes: *Java* (*Apache SOAP* e *Apache Axis*), *Perl* (*SOAP:Lite version 0.52*), *Ruby* (*SOAP4R*), e *C#* no *MS Visual Studio .NET*. O *Google* adverte que a *API* pode ser usada com qualquer outro tipo de linguagem, desde que essa tenha suporte a *web services* (*GOOGLE, 2006*).

<sup>10</sup> *UTF-8* – Do Inglês *8-bit Unicode Transformation Format*. É um tipo de codificação *Unicode* de comprimento variável criado por Ken Thompson e Rob Pike. Pode representar qualquer caractere universal padrão do *Unicode*, sendo também compatível com o *ASCII*. Por esta razão, está lentamente sendo adotado como tipo de codificação padrão para e-mail, páginas *web*, e outros locais onde os caracteres são armazenados. Disponível em <http://www.unicode.org/>.

Para realizar algum tipo de pesquisa através da *API* do *Google* o desenvolvedor deve submeter uma pesquisa ao *Google* com os parâmetros anteriormente citados. O *Google* por sua vez devolve um documento (no formato *XML*<sup>11</sup>) contendo algumas *tags*, conforme tabela 4.3.

Tags	Descrição
<i>documentFiltering</i>	Um valor booleano que indica se um filtro foi utilizado na pesquisa. Será verdadeiro somente quando: (a) Foi requisitado um filtro e (b) o filtro realmente ocorreu.
<i>searchComments</i>	Um texto contendo as palavras mais comuns que foram removidas automaticamente da pesquisa. Estas palavras são conhecidas como “ <i>stop-words</i> ”. Ex.: “e”, “como”, onde, etc.
<i>estimatedTotalResultsCount</i>	Valor estimado do total de resultados para a pesquisa realizada.
<i>estimateIsExact</i>	Um valor booleano indicando se o valor estimado corresponde ao valor real de resultados.
<i>resultElements</i>	Uma lista de itens que corresponde aos resultados da pesquisa.
<i>searchQuery</i>	Onde estarão representados os valores informados em “q”, conforme tabela 4.1.
<i>startIndex</i>	Indica o índice do primeiro resultado da pesquisa
<i>endIndex</i>	Indica o índice do ultimo resultado da pesquisa.
<i>searchTips</i>	Um texto contendo algumas informações de como utilizar o <i>Google</i> .
<i>directoryCategories</i>	Uma lista de itens contendo as categorias de diretórios. Corresponde ao diretório <i>ODP</i> <sup>12</sup> da pesquisa.
<i>searchTime</i>	O tempo em segundos em que a pesquisa foi realizada.

<sup>11</sup> *XML* – Do Inglês *Extensible Markup Language* (Linguagem de Marcação Expansível). É um subtipo de *SGML* (acrônimo de *Standard Generalized Markup Language*, ou *Linguagem Padronizada de Marcação Genérica*) capaz de descrever diversos tipos de dados. Seu propósito principal é a facilidade de compartilhamento de informações através da Internet. Disponível em <http://www.w3.org/XML/>.

<sup>12</sup> *ODP* – Do inglês *Open Directory Project* (Projeto de Diretório Aberto). O *ODP* é o mais amplo e abrangente diretório da *web* editado por humanos. Ele é construído e mantido por uma vasta comunidade global de editores voluntários. Disponível em <http://dmoz.org/>.



<i>summary</i>	Se a pesquisa retornou diretórios <i>ODP</i> , o resumo deles é mostrado aqui.
<i>URL</i>	A <i>URL</i> do resultado da pesquisa.
<i>snippet</i>	Mostra um resumo do local onde os termos de pesquisa estão presentes, dentro da <i>URL</i> .
<i>title</i>	Título do resultado da pesquisa. Geralmente o nome da página.
<i>cachedSize</i>	Texto contendo o tamanho do cachê (inteiro + k). Indica se o cachê da <i>URL</i> está disponível. O tamanho é mostrado em <i>kilobytes</i> .
<i>relatedInformationPresent</i>	Booleano indicando se o termo da pesquisa relacionada é suportado pela <i>URL</i>
<i>hostName</i>	Nome do <i>host</i> .
<i>directoryCategory</i>	Categoria do diretório pesquisado.
<i>directoryTitle</i>	Nome do diretório.
<i>fullViewableName</i>	Texto contendo o nome do diretório <i>ODP</i> , para a corrente categoria.
<i>specialEncoding</i>	Especifica o esquema de codificação da informação do diretório.

**Tabela 4.3 – Tags retornadas pelo Google (GOOGLE, 2006)**

Para maiores informações sobre a utilização das *API's*, o *Google* disponibiliza para os desenvolvedores diversas formas de suporte, como *FAQ's*, exemplos e alguns grupos de pesquisa sobre o assunto.

#### **4.2 Yahoo API**

O mecanismo de busca *Yahoo* também disponibiliza sua vasta base de dados através de suas *API's*. Através destes serviços desenvolvedores podem ter acesso a diversas funcionalidades disponibilizadas pelo mecanismo. As aplicações desenvolvidas devem ser reportadas ao *Yahoo*, como uma forma de manter o nível dos serviços prestados e também devem obedecer a política de uso do mecanismo (*YAHOO*, 2006).

A tabela 4.4 mostra a definição dos parâmetros que devem ser enviados para fazer uma requisição de consulta através da API do Yahoo.

Parâmetro	Valor	Descrição
<i>appid</i>	<i>string</i> (obrigatório)	O identificador da aplicação.
<i>query</i>	<i>string</i> (obrigatório)	O que o usuário deseja pesquisar.
<i>region</i>	<i>string</i> : padrão = <i>us</i>	A região em que se deseja pesquisar. Ex.: Brasil = Br, Argentina = ar, etc.
<i>type</i>	<i>all</i> (padrão), <i>any</i> , or <i>phrase</i>	Tipo de consulta: <ul style="list-style-type: none"> <li>• <i>All</i> – retorna os resultados baseados em todos os termos da pesquisa;</li> <li>• <i>Any</i> – Qualquer resultado, desde que tenha um ou mais termos da pesquisa;</li> <li>• <i>Phrase</i> – Retorna resultados contendo os termos pesquisados em uma única frase.</li> </ul>
<i>results</i>	<i>integer</i> : padrão 10, max 100	O número de resultados que a pesquisa deve retornar.
<i>start</i>	<i>integer</i> : padrão 1	A posição inicial do resultado. A última posição ( <i>start</i> + <i>results</i> – 1) não deve ser superior a 1000.
<i>format</i>	<i>any</i> (padrão), <i>html</i> , <i>mword</i> , <i>pdf</i> , <i>ppt</i> , <i>rss</i> , <i>txt</i> , <i>xls</i>	Especifica que tipo de arquivo está sendo procurado.
<i>adult_ok</i>	<i>no value</i> (padrão), or 1	Especifica se a pesquisa deve retornar conteúdo adulto. O número 1 habilita a pesquisa nesse tipo de site.
<i>similar_ok</i>	<i>no value</i> (padrão), or 1	Especifica se deverão ser omitidos os resultados similares.
<i>language</i>	<i>string</i> : padrão <i>no value</i> (todas)	Linguagem em que será feita a pesquisa.
<i>country</i>	<i>string</i> : padrão <i>no value</i>	O país para restringir a busca somente em sites alocados em cada país.

<i>site</i>	<i>string</i> : padrão <i>no value</i>	Um domínio para restringir uma pesquisa.
<i>subscription</i>	<i>string</i> : padrão <i>no value</i>	Qualquer descrição para um conteúdo Premium que também deve ser consultado. É possível submeter várias descrições.
<i>license</i>	<i>any</i> (padrão), <i>cc_any</i> , <i>cc_commercial</i> , <i>cc_modifiable</i>	<i>Creative Commons license</i> (site da <i>web</i> : <a href="http://creativecommons.org/">http://creativecommons.org/</a> ): licença para utilizar alguns termos.
<i>output</i>	<i>string</i> : <i>xml</i> (padrão), <i>json</i> , <i>php</i>	Formato de retorno da consulta.
<i>callback</i>	<i>String</i>	Parâmetro válido somente quando o formato de saída for igual a <i>json</i> ( <i>JavaScript Object Notation</i> ).

**Tabela 4.4 – Parâmetros de Pesquisa do Yahoo (YAHOO, 2006)**

Assim como o *Google*, o *Yahoo* também limita a utilização da sua *API*. O limite para consultas por *IP* de usuário é de 5000 requisições por dia, sendo que para cada consulta realizada serão permitidos no mínimo 10 e no máximo 100 resultados. Desta forma se o usuário por algum motivo trocar o *IP*, terá o seu limite restabelecido. Outra forma de aumentar o limite, segundo o *Yahoo* (2006), é cadastrar mais de um “*appid*” (*Application ID*) no site. Este *appid* é responsável por identificar a aplicação do usuário no *Yahoo*.

A *API* de busca do *Yahoo* contém bibliotecas que suportam o desenvolvimento nas seguintes linguagens: *Perl*, *Python*, *PHP*, *Java*, *JavaScript* e *Flash*.

Através do preenchimento de um formulário fornecido pelo *Yahoo* o desenvolvedor pode solicitar alterações nas *API*'s disponibilizadas, essas solicitações serão analisadas pelo *Yahoo* e poderão ser realizadas num momento oportuno (YAHOO, 2006).

Para realizar algum tipo de pesquisa através da *API* do *Yahoo*, o desenvolvedor deve submeter uma pesquisa ao *Yahoo* com os parâmetros anteriormente citados. O *Yahoo* por sua vez devolve um documento (no formato *XML*) contendo algumas *tags*, conforme Tabela 4.5.

<i>Tags</i>	Descrição
<i>ResultSet</i>	Contem todos os resultados da pesquisa. Possui os seguintes atributos: <ul style="list-style-type: none"> <li>• <b><i>totalResultsAvailable</i></b>: O número de resultados encontrados na base de dados do mecanismo, para a consulta realizada;</li> <li>• <b><i>totalResultsReturned</i></b>: O número de resultados encontrados e retornados. Este valor pode ser menor que o número de resultados requisitados;</li> <li>• <b><i>firstResultPosition</i></b>: A posição do primeiro resultado de toda a pesquisa.</li> </ul>
<i>Result</i>	Contem cada resposta individual da consulta realizada.
<i>Title</i>	O título da página encontrada na consulta.
<i>Summary</i>	Resumo associado com a página encontrada na consulta.
<i>URL</i>	Endereço da página encontrada na consulta.
<i>ClickURL</i>	O endereço que esta apontando para a página.
<i>MimeType</i>	O tipo <i>MIME</i> <sup>13</sup> da página.
<i>ModificationDate</i>	A data da ultima modificação da página.
<i>Cache</i>	O endereço do cachê e o seu tamanho em <i>bytes</i> .

**Tabela 4.5 – Tags retornadas pelo Yahoo (YAHOO, 2006)**

O *Yahoo* tem diversas formas de suporte para o desenvolvedor que queira utilizar suas *API's* para desenvolvimento de aplicações, entre elas destaca-se o *FAQ* e as listas de e-mail para suporte.

### **4.3 Technorati API**

A *Technorati* possui um dos maiores acervos de *blogs* da *Internet*, contendo 90 milhões de *blogs* cadastrados no seu banco de dados. Segundo dados da própria empresa, por dia são cadastrados mais de 175 mil *blogs*, com um total de 1,6 milhões de *posts* diários.

<sup>13</sup> MIME - Do Inglês *Multipurpose Internet Mail Extensions* (Extensões Multifunção para Mensagens de Internet). É uma norma da internet para o formato das mensagens de correio eletrônico. A grande maioria das mensagens de correio eletrônico são trocadas usando o protocolo *SMTP* e usam o formato *MIME*. As mensagens na *Internet* têm uma associação tão estreita aos padrões *SMTP* e *MIME* que algumas vezes são chamadas de mensagens *SMTP/MIME*. Disponível em <http://www.ietf.org/>.

Através da sua *API* de busca o mecanismo disponibiliza toda esta vasta quantidade de *blogs* para desenvolvedores criarem suas próprias aplicações (*TECHNORATI*, 2006). A *Technorati* incentiva o desenvolvimento de aplicações que utilizem sua *API's*, para tornar seu mecanismo mais conhecido na comunidade de *software*.

Assim como as *API's* anteriores, para se ter acesso às funcionalidades o desenvolvedor precisa criar uma conta no site e solicitar uma chave para utilizar as *API's* em suas aplicações.

A tabela 4.6 traz os parâmetros necessários para se realizar uma pesquisa através da *API* de buscas da *Technorati*.

Parâmetros	Descrição
<i>key</i>	A chave para utilizar a <i>API</i> (obrigatório).
<i>query</i>	O que o usuário deseja pesquisar (obrigatório).
<i>format</i>	Formato de retorno da consulta. <i>XML</i> ou <i>RSS</i> <sup>14</sup> .
<i>language</i>	A região em que se deseja pesquisar. Ex.: Brasil = Br, Argentina = ar, etc.
<i>authority</i>	Uma forma de fazer um ranking dos <i>blogs</i> mais visitados contidos no banco de dados. Possui as seguintes opções: <ul style="list-style-type: none"> <li>• <i>n</i> – Todos os resultados;</li> <li>• <i>a1</i> – <i>Blogs</i> que contém pelo menos um link apontando pra eles;</li> <li>• <i>a4</i> – <i>Blogs</i> que contém muitos links apontado para eles;</li> <li>• <i>a7</i> – <i>Blogs</i> que contém milhares de links apontando para eles.</li> </ul>
<i>start</i>	Onde quer se começar a ver os resultados, obrigatoriamente deve ser maior que zero.
<i>limit</i>	Limite de <i>links</i> retornados na pesquisa. Mínimo 1 e máximo 100.
<i>Claim</i>	Informações do proprietário do <i>blog</i> pesquisado. '1' significa que

<sup>14</sup> *RSS* - Arquivo em *XML* que serve para compartilhar informações e notícias entre sites e entre sites e programas. Disponível em <http://www.rssfeeds.com.br>.

	deve mostrar e '0' não mostrar.
--	---------------------------------

**Tabela 4.6 – Parâmetros de Pesquisa do *Technorati* (TECHNORATI, 2006)**

Conforme tabela acima a *API* possui algumas limitações, como por exemplo, o número de links que podem ser retornados, 100 links por pesquisa. Não existe um limite de consultas que podem ser feitas por dia ou por *IP* como as *API's* do *Google* e *Yahoo*.

A *API* da *Technorati* tem suporte ao desenvolvimento de aplicações através das seguintes linguagens: *C#, Java, Perl, PHP, Python, REBOL, Ruby, e Visual Basic*.

Para realizar algum tipo de pesquisa através da *API* do *Technorati* o desenvolvedor deve submeter uma pesquisa ao mecanismo com os parâmetros anteriormente citados. O *Technorati* por sua vez devolve um documento (no formato *XML* ou *RSS*) contendo algumas *tags*, conforme tabela 4.7.

<i>Tags</i>	Descrição
<i>query</i>	Palavras pesquisadas.
<i>Querycount</i>	Número de links encontrados
<i>Querytime</i>	Tempo em que a pesquisa retornou os resultados.
<i>Rankingstart</i>	Posição em que a pesquisa irá começar.
<i>name</i>	Nome do <i>blog</i> onde foram encontradas as palavras informadas pelo usuário.
<i>url</i>	Endereço da página.
<i>rssurl</i>	Endereço do <i>RSS</i> da página.
<i>atomurl</i>	Endereço do <i>Atom</i> <sup>15</sup>
<i>inboundblogs</i>	<i>Blogs</i> contidos no resultado
<i>inboundlinks</i>	Links contidos no resultado.
<i>lastupdate</i>	Última atualização do <i>Blog</i> .

<sup>15</sup> *Atom* - Nome dado a um estilo baseado em conteúdo *XML* e meta data, ou seja, é um protocolo ao nível da aplicação para publicar e editar fontes *web* que são periodicamente atualizadas, como por exemplo *Blogs*. Os *feeds* devem ser formados em formato *XML* e são identificados como *application/atom+xml media type*. Disponível em <http://www.atomenabled.org/>.

<i>title</i>	Título da entrada do <i>Blog</i> .
<i>excerpt</i>	Resumo contendo as palavras pesquisadas destacadas.
<i>created</i>	Quando foi criada a entrada.
<i>permalink</i>	Endereço para a entrada do <i>blog</i> .

**Tabela 4.7 – Tags retornadas pelo *Technorati* (TECHNORATI, 2006)**

#### 4.4 Comparativo entre as API's de Mecanismos de Busca

Após estudo das principais características das três API's de mecanismos de busca anteriormente citados, é possível apontar algumas diferenças básicas, como por exemplo, as limitações sobre as pesquisas. A API do *Google* limita suas consultas em apenas dez por dia, enquanto através da API do *Yahoo* é possível realizar até cinco mil consultas diárias. Neste ponto a API do *Technorati* supera as duas anteriores, pois não limita o número de consultas. Uma outra limitação que ocorre na API do *Google*, e que não ocorre nas demais, diz respeito à quantidade de termos que podem ser informados por consulta, no *Google* são somente dez termos. Outra limitação que pode ser destacada é referente à quantidade de resultados retornados por pesquisa realizada, enquanto a API do *Google* pode retornar no máximo dez resultados, as API's do *Yahoo* e *Technorati* retornam até cem resultados. A tabela 4.8 mostra este comparativo.

	Número de Consultas Diárias	Número de Termos por Consulta	Número de Resultados Retornados
<i>Google</i>	Até 10	Até 10	Até 10
<i>Yahoo</i>	Até 5.000	Sem Restrição	Até 100
<i>Technorati</i>	Sem Restrição	Sem Restrição	Até 100

**Tabela 4.8 – Comparativo entre os Mecanismos de Busca**

Apesar da API do *Technorati* ter limitações menores em relação às demais, sua utilização é mais complexa, pois não disponibiliza métodos muito claros para possibilitar a construção de aplicações, ou seja, somente após muitos testes é possível começar a compreender seu funcionamento. As API's do *Google* e *Yahoo* disponibilizam uma documentação muito clara sobre todas as funcionalidades que podem ser implementadas a partir das suas API's, inclusive mostrando exemplos das mesmas, desta forma facilitando a utilização das mesmas para o desenvolvimento de aplicações.

## 5 PROCESSAMENTO DA LINGUAGEM NATURAL (PLN)

Este capítulo versa sobre alguns conceitos a respeito de PLN, um pouco do seu histórico e principalmente o estudo de alguns tipos de ferramentas que utilizam seus recursos. Este estudo é utilizado para dar uma base para a construção do *framework* proposto. As ferramentas estudadas estão classificadas nas seguintes categorias:

- *Tagger, Morphological Analyzer* - Analisadores morfológicos e etiquetadores;
- *Stemmer* – Radicalizadores de palavras;
- *Parsers* - Analisadores sintáticos;
- *Corpus Tools* - Ferramentas para Manipulação de *Corpus*.

### 5.1 Conceitos de PLN

“A PLN visa tratar a linguagem natural com o objetivo de desenvolver teorias computacionais da linguagem, usando noções de algoritmos e estrutura de dados provenientes da ciência da computação” (ABRAHÃO, 1997, p. 10).

Para Chaves (2003a) o PLN é o conjunto de métodos formais para analisar textos e gerar frases escritas em um idioma humano.

Segundo Vieira e Lima (2001) PLN envolve a construção de programas capazes de interpretar e/ou gerar informação fornecida em linguagem natural. Para o processamento da língua natural, são necessários vários subsistemas para dar conta dos diferentes aspectos da língua: sons, palavras, sentenças e discurso nos níveis estruturais, de significado e de uso.



Para Nunes et al. (1999) PLN compreende a criação de programas computacionais “inteligentes”, até mesmo capazes de “compreender” as línguas e, por meio delas, simular uma interação verbal com o usuário.

“O PLN é uma forma artificial desenvolvida com a finalidade de aprender através da linguagem natural (LN), artifícios que possam ser utilizados com recursos computacionais, para desenvolver programas capazes de tratar os problemas relacionados à LN” (VINHAES, 2005, p.20).

Segundo Jackson e Moulinier (2002), o termo Processamento de Linguagem Natural (PLN) é normalmente usado para descrever a função de *softwares* ou componentes de *hardware* em um sistema computacional que analisam ou sintetizam linguagem falada ou escrita. O termo “natural” tem a função de diferenciar a fala e a escrita humana de linguagens mais formais, tais como notações matemáticas ou lógicas, ou ainda linguagens de computador, tais como Java, LISP e C++.

## 5.2 Histórico do PLN

Na década de 50 a comunidade científica começava com os primeiros estudos sobre PLN. Segundo Nunes et al. (1999) as primeiras investigações institucionalizadas sobre o Processamento da Linguagem Natural começaram a ser desenvolvidas no início da década de 50, depois da distribuição de 200 cópias de uma carta, conhecida como *Weaver Memorandum*, escrita por Warren Weaver, então vice-presidente da Fundação Rockefeller e exímio conhecedor dos trabalhos sobre criptografia computacional. Nessa carta, divulgada em 1949, Weaver convidava universidades e empresas, interessados potenciais, para desenvolver projetos sobre um novo campo de pesquisa que ficou conhecido como “tradução automática”, “tradução mecanizada” ou simplesmente MT (abreviação do inglês “*Machine Translation*”).

Em 1960 foram desenvolvidas as primeiras formalizações do significado em termos de redes semânticas, os primeiros tratamentos computacionais das gramáticas livres de contexto e a criação dos primeiros analisadores. Na década de 70 veio à consolidação dos estudos do PLN com a implementação de parcelas das primeiras gramáticas e analisadores

sintáticos e a busca de formalização de fatores pragmáticos e discursivos. Em 1980 começava a sofisticação dos sistemas com o desenvolvimento de novas teorias lingüísticas motivadas pelo estudo do PLN. Nos anos 90 surgiam os sistemas baseados em representação do conhecimento com o desenvolvimento de projetos de sistemas de PLN complexos que buscavam a integração dos vários tipos de conhecimentos lingüísticos e extralingüísticos e das estratégias de inferência envolvidas nos processos de produção, manipulação e interpretação de objetos lingüísticos (NUNES et al., 1999).

### **5.3 *Tagger, Morphological Analyser* (Etiquetador, Analisador Morfológico)**

A morfologia além de estudar a estrutura das palavras também se preocupa com a classificação das palavras em categorias, ou, conforme um termo muito conhecido na área de PLN, as palavras são classificadas em partes do discurso (*part-of-speech*, ou *POS*). Essas categorias podem ser divididas em substantivos, verbos, adjetivos, preposições e advérbios. Ainda pode-se dividir as categorias em classes abertas ou fechadas. Abertas são constituídas por categorias que abrangem um grande número de palavras e podem, ainda, abrigar o surgimento de palavras novas. Nessa classe encontram – se os substantivos, verbos e adjetivos. Classes fechadas são aquelas que possuem funções gramaticais bem definidas, como por exemplo, artigos, demonstrativos, quantificadores, conjunções e preposições (VIEIRA; LIMA, 2001; CHAVES, 2003a).

Sendo assim a função de um etiquetador morfológico é associar a cada palavra de uma sentença uma determinada etiqueta, que corresponda a sua categoria morfológica ou gramatical (MENEZES; NETO, 2002).

Para a atribuição das etiquetas morfológicas o etiquetador realiza as seguintes tarefas: associa etiquetas possíveis às palavras da sentença, avalia a probabilidade de cada etiqueta e define a etiqueta mais provável para cada palavra da sentença (examina o contexto) (KINOSHITA et al., 2005).

Uma das grandes dificuldades existentes na tarefa da classificação morfológica encontra-se em sua susceptibilidade à ambigüidade. Um etiquetador morfológico robusto deve levar em conta não apenas as informações lexicais da palavra a ser anotada, mas também

informações a respeito do contexto em que está inserida (MENEZES; NETO, 2002; VIEIRA; LIMA, 2001; ROCHE; SCHABES, 1995).

#### **5.4 Stemmer (Radicalizador de Palavras)**

Conflação é o ato de fusão ou combinação para realizar o processo de igualar variantes morfológicas de termos. A conflação pode ser feita de duas formas, uma é a manual, utilizando algum tipo de expressão regular, ou ainda automática, através do uso de programas chamados *stemmers*. Para reduzir as variantes ou variações de uma palavra para uma forma única são utilizados algoritmos de conflação. Existem dois métodos distintos de conflação para se obter tal redução: *Stemming* e redução à forma canônica ou lematização (CHAVES, 2003b; GONZALEZ, 2005)

*Stemming* consiste no processo de combinar as formas diferentes de uma palavra em uma representação comum, chamada de radical (*stem*). Ele não é necessariamente igual a raiz lingüística, porém permite tratar diferentes variações de uma palavra da mesma forma. Por exemplo, “conector” e “conectores” são essencialmente iguais, mas se sofrerem a redução por *stemming* serão tratadas como palavras diferentes (DIAS; MALHEIROS, 2006; CHAVES, 2003b). Como pode ser observado a seguir.

No processo de *stemming* podem ser detectados dois tipos de erros, conhecidos como *overstemming* e *understemming*. *Overstemming* ocorre quando a parte removida da palavra não é um sufixo, mas parte do *stem*. Desta forma palavras não relacionadas podem ser combinadas. Por exemplo, a palavra “confortável” após ter sido processada por um radicalizador, é transformada no *stem* “confor”. Como pode ser observado foi removido parte do radical correto, a saber, “confort”. *Understemming* ocorre quando um sufixo não é removido completamente. Por exemplo, quando a palavra “referência” é transformada incorretamente no *stem* “referênc”, quando deveria ter sido transformada no *stem* “refer” (DIAS; MALHEIROS, 2005; CHAVES, 2003b).

Outro problema que se pode observar no processo de *stemming* é que ele não tem sucesso com termos onde a flexão é raramente usada ou inexistente (por exemplo, nomes próprios) (GONZALEZ, 2005).

A forma canônica ou lema é a forma pela qual constam as palavras na nomenclatura de um dicionário usual: masculino singular, para os nomes e adjetivos (pato, gordo), feminino singular, para nomes e adjetivos exclusivamente femininos (mesa, grávida) e infinitivo para os verbos (lavar, comer, fugir). Por exemplo, seguindo as regras acima, a forma canônica de gatos e gatas é gato (BAPTISTA, 1994; NUNES et al., 1996).

Um problema que pode ser detectado no método de lematização é que algumas palavras da mesma família morfológica podem ser classificadas como sendo diferentes. Por exemplo, a palavra “caminhei” ficaria igual a “caminhar”, enquanto a palavra da mesma família “caminhada” ficaria igual a “caminhada” (GONZALEZ, 2005).

### **5.5 Parser (Analisador Sintático)**

O processamento automático da sentença com o objetivo do reconhecimento de sua estrutura sintática é chamado de *parsing*. Por extensão, a ferramenta que executa esse conjunto de procedimentos (que permite assinalar funções sintáticas a cada um dos itens lexicais que compõem a sentença) é referenciada como *parser* (NUNES, 1999).

Através da utilização de um *parser* é possível construir uma árvore de derivação, que explicita as relações entre as palavras que formam uma sentença. O analisador sintático faz uso do léxico, que reúne o conjunto de itens lexicais da língua, e de uma gramática, que define as regras de combinação dos itens na formação das frases (sintagmas<sup>16</sup>) (VIEIRA; LIMA, 2001).

Segundo Silva e Chapetta (2003) todo *parser* deve incluir uma gramática que especifique como as sentenças serão interpretadas. Pode-se considerar esta gramática como sendo um espaço de configurações, onde cada configuração representa estágios da combinação dos constituintes da sentença que está sendo interpretada. Transições entre configurações descrevem como os constituintes são combinados para derivar outros constituintes mais extensos.

---

<sup>16</sup> Sintagma é a unidade da análise sintática composta de um núcleo e de outros termos que a ele se juntam, formando uma locução que irá compor a oração. O nome do sintagma depende da classe da palavra que forma seu núcleo, havendo assim sintagma nominal (núcleo nome), sintagma adjetivo (núcleo adjetivo), sintagma preposicional (núcleo preposição) etc (COELHO et al., 2005).

## 5.6 *Corpus Tools*

*Corpus Tools* são ferramentas destinadas à manipulação e criação de corpus para PLN. A seguir serão apresentadas algumas definições sobre o significado de um *corpus*:

“Uma coletânea de textos naturais (*naturally occurring*), escolhidos para caracterizar um estado ou variedade de linguagem” (SINCLAIR, 1991, p. 171).

“[*Corpus* é] um corpo de linguagem natural (autêntica) que pode ser usado como base para pesquisa lingüística” (SINCLAIR, 1991, p. 171).

“*Corpus* é uma coletânea de porções de linguagem que são selecionadas e organizadas de acordo com critérios lingüísticos explícitos, a fim de serem usadas como uma amostra da linguagem” (PERCY et al., 1996, p. 4).

“Um corpo de material lingüístico que existe em formato eletrônico e que pode ser processado por computador para vários propósitos” (LEECH; MACENERY, 1997, p. 1).

“*Corpus* de material lingüístico natural (textos inteiros, amostra de textos, ou às vezes somente sentenças desconexas), que são armazenadas em formato legível por máquina” (LEECH, 1991, p. 115-116).

“Uma coletânea grande e criteriosa de textos naturais (BIBER et al., 1998 , p. 4)”.

[...] um conjunto de dados lingüísticos (pertencentes ao uso oral ou escrito da língua, ou a ambos), sistematizados segundo determinados critérios, suficientemente extensos em amplitude e profundidade, de maneira que sejam representativos da totalidade do uso lingüístico ou de algum de seus âmbitos, dispostos de tal modo que possam ser processados por computador, com a finalidade de propiciar resultados vários e úteis para a descrição e análise [...] (SANCHEZ, 1995, p. 8-9).

“Entende-se por *corpus* um conjunto de textos utilizados, por um sistema de PLN, para treinamento e teste de modelos estatísticos ou para avaliação de componentes de sistemas de PLN” (SANT’ANNA, 2000, p. 36).

### 5.6.1 Tipos de *Corpus*

De acordo com o seu tipo, um *corpus* pode conter textos de gêneros diferentes ou de mesmo gênero. Do mesmo modo, um *corpus* pode ser classificado como marcado, caso possua informações sintáticas ou semânticas para cada item lexical, ou não-marcado, caso essas informações não estejam disponíveis (SANT'ANNA, 2000).

Um *corpus*, de acordo com o uso, pode ser classificado como (SANT'ANNA, 2000):

- *Corpus* de treinamento: usado nos casos em que o programa extrai informações estatísticas de um conjunto de tamanho adequado de textos. Serve para a retirada desses padrões estatísticos;
- *Corpus* de teste: usado para teste do sistema, após a etapa de treinamento;
- *Corpus* de avaliação: usado quando componentes de um sistema precisam ser avaliados.

### 5.6.2 Pré-Requisitos para Formação de um *Corpus* Computadorizado

Segundo Sardinha (2000) os quatro pré-requisitos para formação de um *corpus* computadorizado são:

- Em primeiro lugar, o *corpus* deve ser composto de textos autênticos, em linguagem natural. Desta forma, os textos não podem ter sido produzidos com o propósito de serem alvo de algum tipo de pesquisa lingüística. E não podem ter sido criados em linguagem artificial, tais como linguagem de programação de computadores ou notação matemática.
- Em segundo lugar, quando se fala em autenticidade dos textos, subentendem-se textos escritos por falantes nativos. Tanto assim que, quando este não é o caso, deve-se qualificá-lo, falando-se em *corpora* “de aprendizes” (“*learner corpora*”).

- O terceiro pré-requisito diz que o conteúdo do *corpus* deve ser escolhido criteriosamente. Os princípios da escolha dos textos devem seguir, acima de tudo, as condições de naturalidade e autenticidade. Mas devem também obedecer a um conjunto de regras estabelecidas pelos seus criadores de modo que o *corpus* coletado corresponda às características que se deseja dele. Ou seja, o conteúdo do *corpus* deve ser selecionado a fim de garantir que o *corpus* tenha uma certa característica.

- O quarto pré-requisito é considerado como sendo o mais problemático: Representatividade. Tradicionalmente, um *corpus* é visto como um conjunto representativo de uma variedade lingüística ou mesmo de um idioma. É problemático, pois não se consegue medir o quanto um *corpus* é representativo em relação a outro.

## 6 FERRAMENTAS DE PLN

Este capítulo mostra algumas ferramentas relacionadas com PLN. São apresentadas quatro ferramentas distintas, sendo uma para cada categoria apresentada no capítulo anterior: *Tagger*, *Morphological Analyzer*, *Stemmer*, *Parsers* e *Corpus Tools*.

Neste mesmo capítulo é apresentada uma *API* que implementa alguns algoritmos de PLN, a *API Lucene*.

### 6.1 *WebJspell*

Tradicionalmente, os analisadores morfológicos contêm um dicionário que serve como base (palavras e suas características morfológicas) e um conjunto de regras de formação de novas palavras. Esta idéia é a base do analisador morfológico *Jspell* (SIMÕES; ALMEIDA, 2002).

O *WebJspell* é um serviço desenvolvido pelo Pólo de Braga<sup>17</sup> da Linguateca<sup>18</sup> com objetivo de tornar o analisador morfológico e corretor ortográfico *Jspell* acessível a um maior número de utilizadores. O *Jspell* foi desenvolvido pelo Projeto Natura<sup>19</sup> (LINGUATECA, 2006).

O *WebJspell* e o *Jspell* estão em constante evolução de forma a aumentar a qualidade. É desenvolvido e mantido por Rui Vilela, José João Almeida e Alberto Simões (LINGUATECA, 2006).

---

<sup>17</sup> Pólo de Braga - O Pólo de Braga da Linguateca é um dos 5 pólos da Linguateca. O Pólo de Braga é composto pelo investigador Rui Vilela e o aluno de doutorado Alberto Simões, orientados por José João Dias de Almeida e Diana Santos. O Pólo está sediado no Departamento de Informática da Universidade do Minho. Nasceu da colaboração entre a Linguateca e o Projeto Natura.

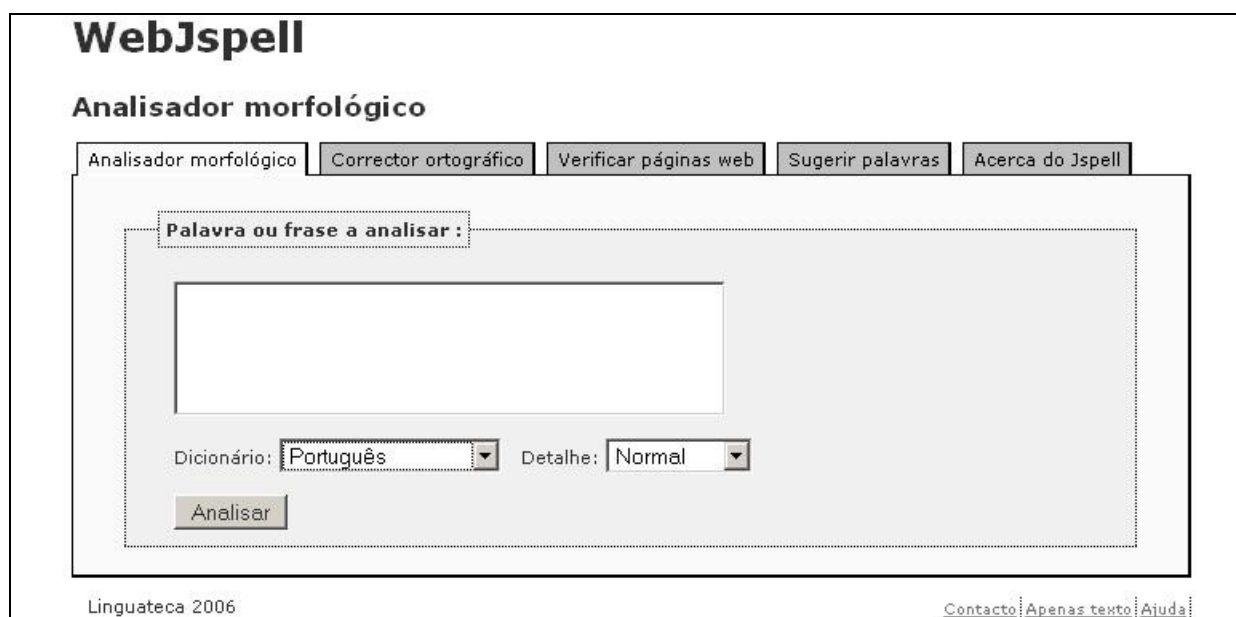
<sup>18</sup> Linguateca – É formalmente um Centro de Recursos distribuído para a Língua Portuguesa. Resumidamente, a missão da Linguateca é promover a existência, discussão, avaliação, distribuição e manutenção de recursos relacionados com o processamento lingüístico do Português.

<sup>19</sup> Projeto Natura - O Projeto Natura contempla processamento de Linguagem Natural (PLN) com especial ênfase no português.



O *Jspell* foi desenvolvido com base no corretor ortográfico *open-source Inspell*. Embora o *Inspell* não seja um analisador morfológico ele já incluía a possibilidade de definição e uso de regras morfológicas elementares. Em vez de se utilizar dicionários de extensão, define-se um conjunto de regras morfológicas (tabela de afixos que permite, a partir de uma palavra, obter várias flexões e derivações) e associa-se a cada palavra o conjunto de regras morfológicas que lhe são aplicáveis. Esta associação permite que um dicionário seja substancialmente menor do que a listagem de todas as palavras dele obtidas e também permite uma maior riqueza no tratamento de palavras desconhecidas (SIMÕES; ALMEIDA, 2002).

A seguir é apresentada a tela inicial (figura 6.1) do *WebJspell*, onde o usuário tem a opção de escolher qual dicionário quer utilizar (Latim, Português ou Inglês) e qual o nível de detalhe da resposta da ferramenta (Normal ou Resumida) ele quer visualizar. Na parte central da ferramenta encontra-se o local onde deverá ser informada a palavra ou frase que se deseja analisar morfológicamente.



**Figura 6.1 – Tela Inicial do *WebJspell***

A seguir são apresentadas algumas figuras com o resultado da análise morfológica da frase “O Futebol deveria ser a alegria do povo.”. Observa-se que o analisador separa cada parte do texto e analisa separadamente.

O	Adicionar análise	Remover do dicionário	Flexionar	Exemplos de frases	
<ul style="list-style-type: none"> <li>• <b>Categoria:</b> artigo</li> <li>• <b>Número:</b> singular</li> <li>• <b>Género:</b> masculino</li> <li>• <b>Lema:</b> o</li> <li>• <a href="#">Formas flexionadas do lema</a></li> <li>• <a href="#">Exemplos de frases de formas flexionadas do lema</a></li> <li>• <a href="#">Exemplos de frases da palavra com o lema corrente</a></li> </ul>					
<ul style="list-style-type: none"> <li>• <b>Categoria:</b> pronome pessoal</li> <li>• <b>Número:</b> singular</li> <li>• <b>Caso latino:</b> acusativo</li> <li>• <b>Pessoa:</b> terceira</li> <li>• <b>Género:</b> masculino</li> <li>• <b>Lema:</b> o</li> <li>• <a href="#">Formas flexionadas do lema</a></li> <li>• <a href="#">Exemplos de frases de formas flexionadas do lema</a></li> <li>• <a href="#">Exemplos de frases da palavra com o lema corrente</a></li> </ul>					

**Figura 6.2 – Análise Morfológica de “O”**

Futebol	Adicionar análise	Remover do dicionário	Flexionar	Exemplos de frases							
<ul style="list-style-type: none"> <li>• <b>Categoria:</b> substantivo comum</li> <li>• <b>Número:</b> singular</li> <li>• <b>Género:</b> masculino</li> <li>• <b>Lema:</b> futebol</li> <li>• <a href="#">Formas flexionadas do lema</a></li> <li>• <a href="#">Exemplos de frases de formas flexionadas do lema</a></li> <li>• <a href="#">Exemplos de frases da palavra com o lema corrente</a></li> </ul>											
<table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td colspan="3"><b>Formas flexionadas do lema <i>futebol</i>:</b></td> </tr> <tr> <td>futebolista</td> <td>futebol</td> <td>futebolistas</td> </tr> </table>						<b>Formas flexionadas do lema <i>futebol</i>:</b>			futebolista	futebol	futebolistas
<b>Formas flexionadas do lema <i>futebol</i>:</b>											
futebolista	futebol	futebolistas									

**Figura 6.3 – Análise Morfológica de “Futebol”**

deveria	Adicionar análise	Remover do dicionário	Flexionar	Exemplos de frases	
<ul style="list-style-type: none"> <li>• <b>Categoria:</b> verbo</li> <li>• <b>Tempo:</b> condicional</li> <li>• <b>Número:</b> singular</li> <li>• <b>Pessoa:</b> primeira/terceira</li> <li>• <b>Transitividade:</b> transitivo/intransitivo</li> <li>• <b>Lema:</b> dever</li> <li>• <a href="#">Tabela de conjugação verbal</a></li> <li>• <a href="#">Formas flexionadas do lema</a></li> <li>• <a href="#">Exemplos de frases de formas flexionadas do lema</a></li> <li>• <a href="#">Exemplos de frases da palavra com o lema corrente</a></li> </ul>					

**Figura 6.4 – Análise Morfológica de “deveria”**

<ul style="list-style-type: none"> <li>• Nenhuma acção associada</li> </ul>					
<ul style="list-style-type: none"> <li>• <b>Categoria:</b> pontuação</li> </ul>					

**Figura 6.5 – Análise Morfológica de “.”**

Analisando as figuras apresentadas pode-se verificar que alguns tipos morfológicos tem um tipo de apresentação diferente na ferramenta *WebJspell*. Para a categoria verbo nota-se a presença de um link que aponta para a conjugação verbal da palavra.

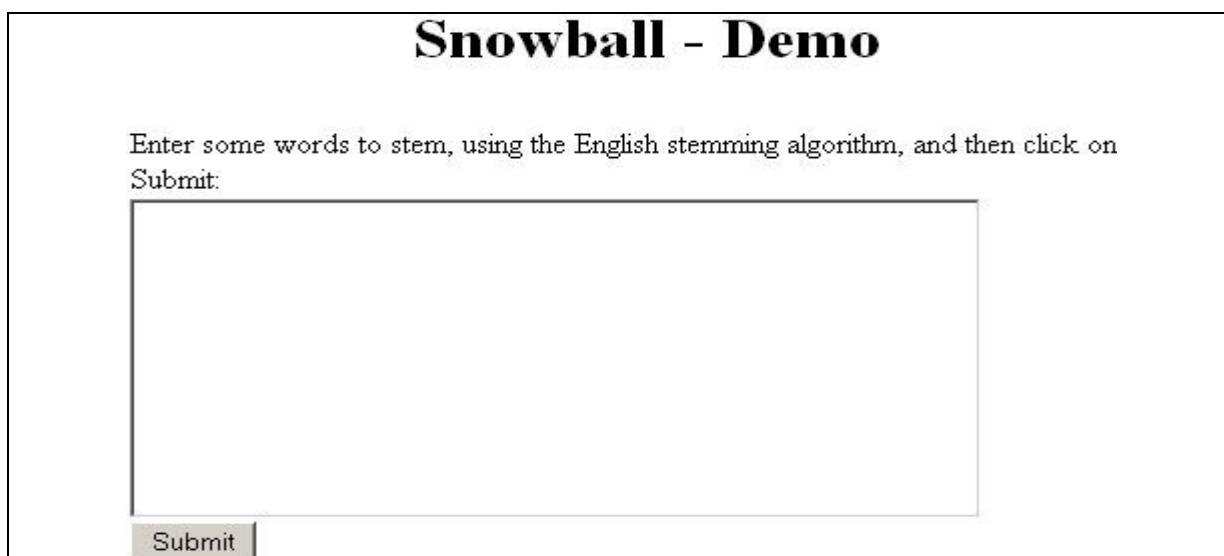
Uma das grandes limitações encontradas em qualquer ferramenta de análise morfológica é a limitação do seu dicionário, pois sempre existirão palavras que não constam nos dicionários armazenados nas ferramentas.

## 6.2 *Snowball*

O *Snowball* (SNOWBALL, 2006) é um *stemmer* que possui módulos em várias línguas, sendo uma delas o português. Através de vocabulários armazenados ele processa as palavras informadas e mostra as palavras na sua forma canônica ou radicalizada.

O *Snowball* é baseado em um algoritmo desenvolvido em uma micro-linguagem de mesmo nome. Os algoritmos utilizados para a criação do *Snowball* foram desenvolvidos por Martin Porter (SNOWBALL, 2006).

A figura 6.6 mostra a ferramenta *Snowball*, nesse *demo* somente é possível fazer a radicalização na língua inglesa. A figura 6.7 mostra o resultado obtido através da ferramenta para a palavra “*Several*”.

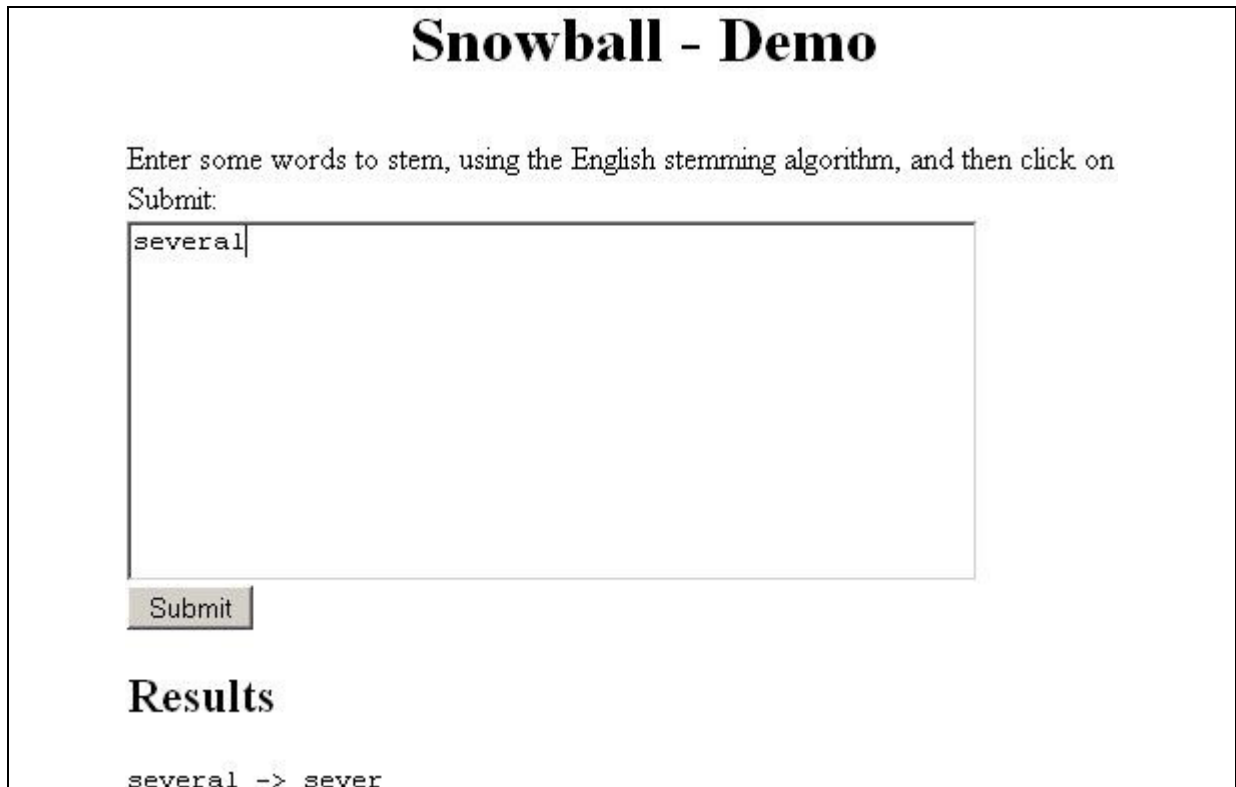


**Snowball - Demo**

Enter some words to stem, using the English stemming algorithm, and then click on Submit:

Submit

**Figura 6.6 – *Snowball***



**Figura 6.7 – Resultado de Análise do Snowball**

Através do resultado pode-se observar que o método utilizado pelo *Snowball* é o *stemming*, onde a palavra é reduzida a seu radical. Observa-se também que a palavra *several* (do inglês vários) não poderia ser radicalizada para *sever* (do inglês cortar algo, romper), este já é um problema encontrado nesta ferramenta, pois a palavra perderia seu significado.

### **6.3 FLIP On-Line**

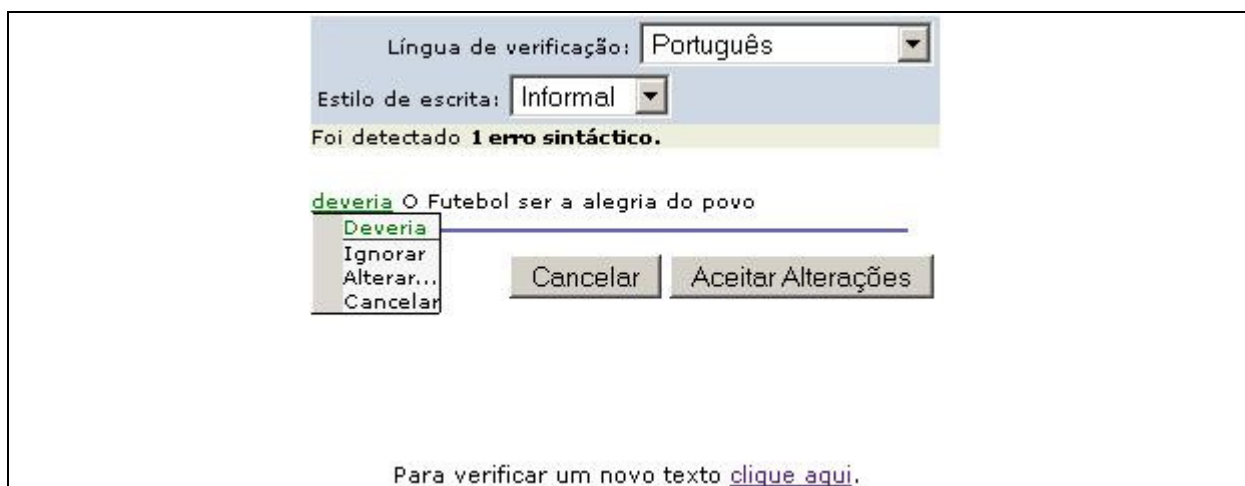
O *Flip on-line* é uma ferramenta na categoria de *parsers* desenvolvida pela empresa Priberam. Segundo *Flip* (2006) a ferramenta foi desenvolvida para facilitar o acesso à tecnologia *Flip*. A Priberam disponibilizou este serviço gratuito *on-line* e também disponibilizou o corretor ortográfico, tanto para português europeu quanto para português do Brasil. A partir da janela disponível na página, o usuário pode digitar as palavras ou as frases (até ao limite de 3000 caracteres) sobre as quais tem dúvidas, selecionar a variedade de português que pretende e visualizar as correções propostas.

A figura a seguir mostra a janela inicial da ferramenta.



**Figura 6.8 – Flip On-Line**

A figura 6.9 mostra o resultado da análise sintática feita sobre uma frase que está sintaticamente incorreta. Observa-se também que existe a opção de fazer análise sobre a ótica do estilo de escrita, como Formal, Corrente e Informal.



**Figura 6.9 – Análise Sintática do Flip On-Line**

O *Flip* apontou um erro sintático nesta frase, o verbo “deveria” deveria estar com a primeira letra maiúscula. Na figura a seguir nota-se que existem dois erros ortográficos que o *parser* identificou.

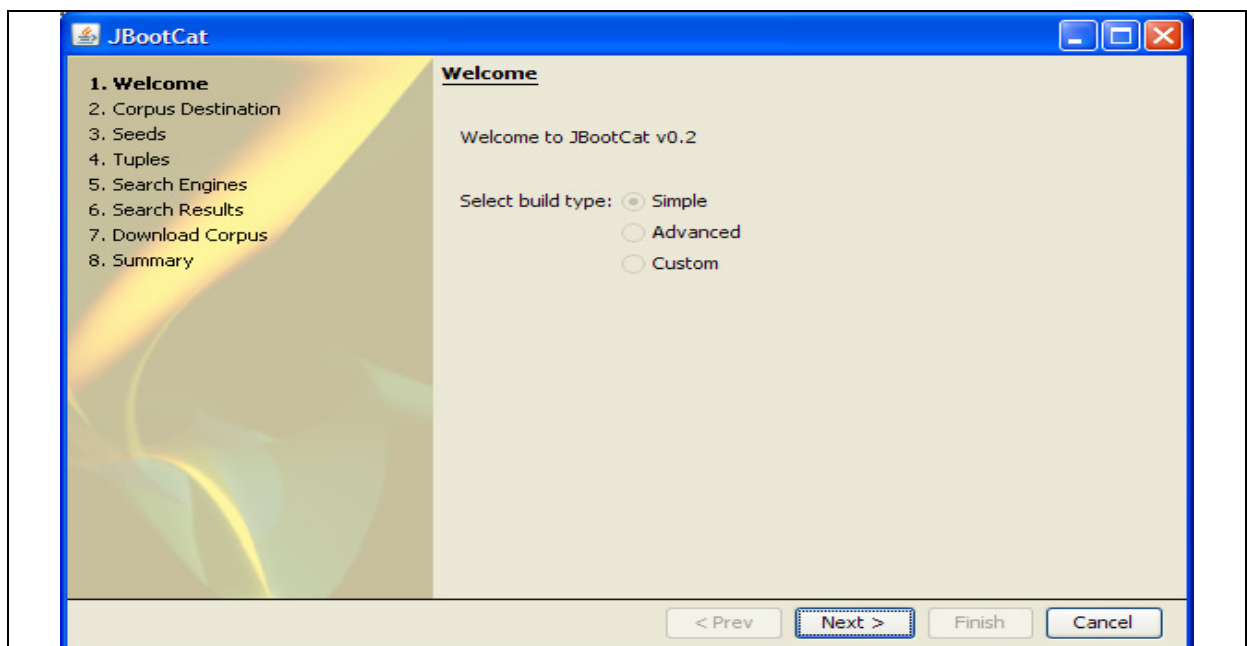


**Figura 6.10 – Análise Sintática do *Flip On-Line***

#### 6.4 *JBootCat*

O *JBootCat* (ROBERTS, 2006) é uma ferramenta que se enquadra na categoria de *Corpus Tools*. Esta ferramenta tem uma proposta semelhante com a do corrente trabalho, pois utiliza a *web* como base de dados para gerar o *corpus*.

A seguir é apresentada a tela inicial da ferramenta que foi desenvolvida utilizando scripts escritos por Marco Baroni.



**Figura 6.11 – Tela inicial do *JBootCat***

Através desta ferramenta o usuário informa o caminho onde será salvo o novo *corpus* que será criado pelo usuário (figura 6.12), as palavras que deseja encontrar (figura 6.13), essas

palavras ainda poderão ser salvas para uma consulta futura. Após a escolha das palavras o usuário deverá gerar as tuplas (poderá escolher qual o número de tuplas e quantas palavras deverão constar em cada uma) para a pesquisa no motor de busca (figura 6.14). O próximo passo é escolher o motor de buscas na *web* (somente o *Google* por enquanto) e inserir a chave de pesquisa (conforme o capítulo sobre *API'S*) (figura 6.15). No passo seguinte o usuário deverá solicitar que a ferramenta efetue a pesquisa na *web*, e com a pesquisa realizada o mesmo deverá selecionar quais os resultados que deverão constar no *corpus* (figura 6.16). Na etapa seguinte o usuário deverá solicitar que o sistema faça o *download* do *corpus* com as informações que ele selecionou (figura 6.17). Na etapa final é apresentado um resumo das ações realizadas pela ferramenta *JBootCat* (figura 6.18).

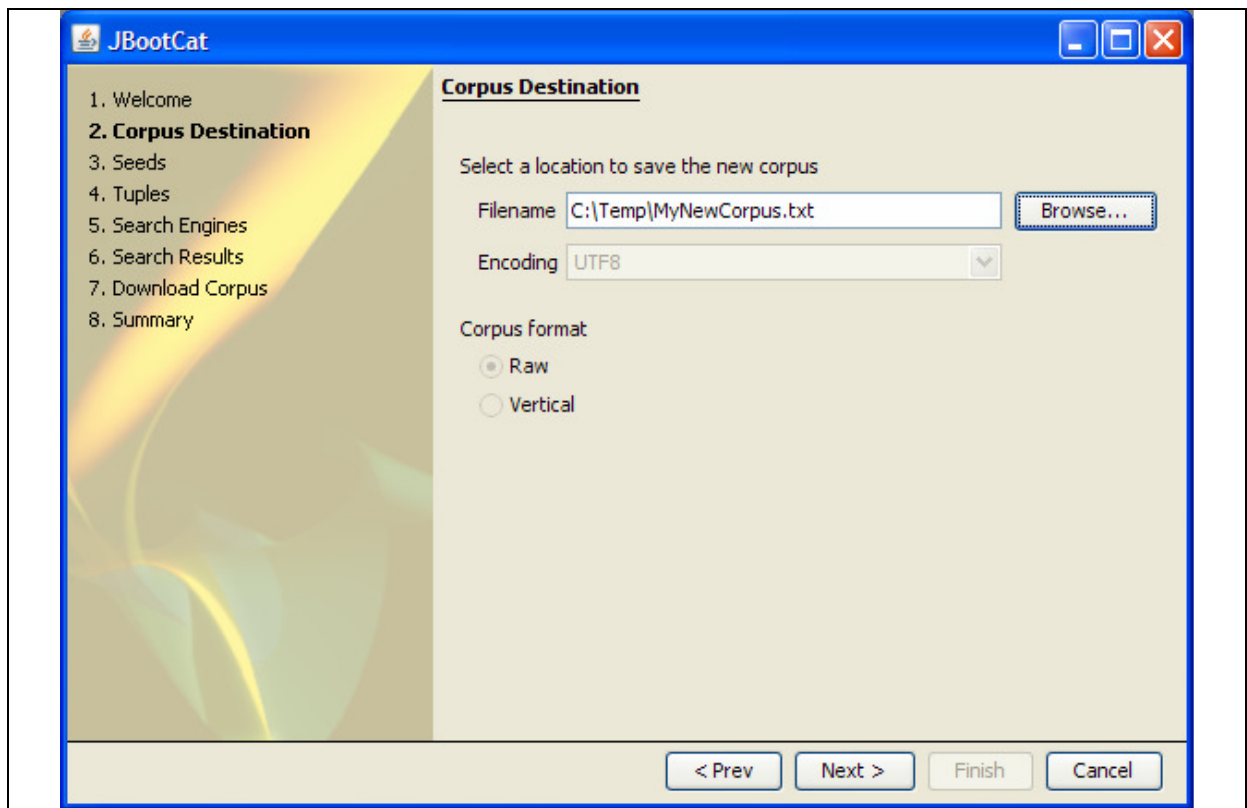


Figura 6.12 – Escolha do local para salvar o novo *corpus*

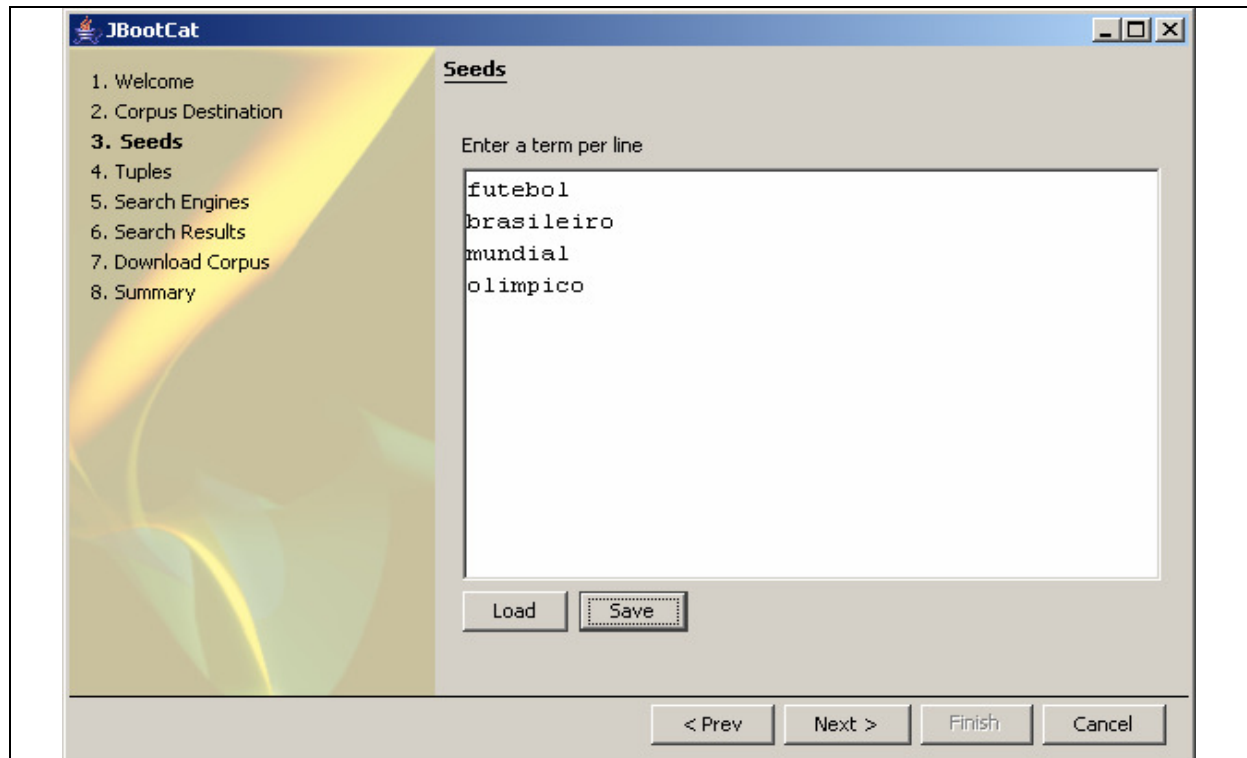


Figura 6.13 – Escolha das palavras

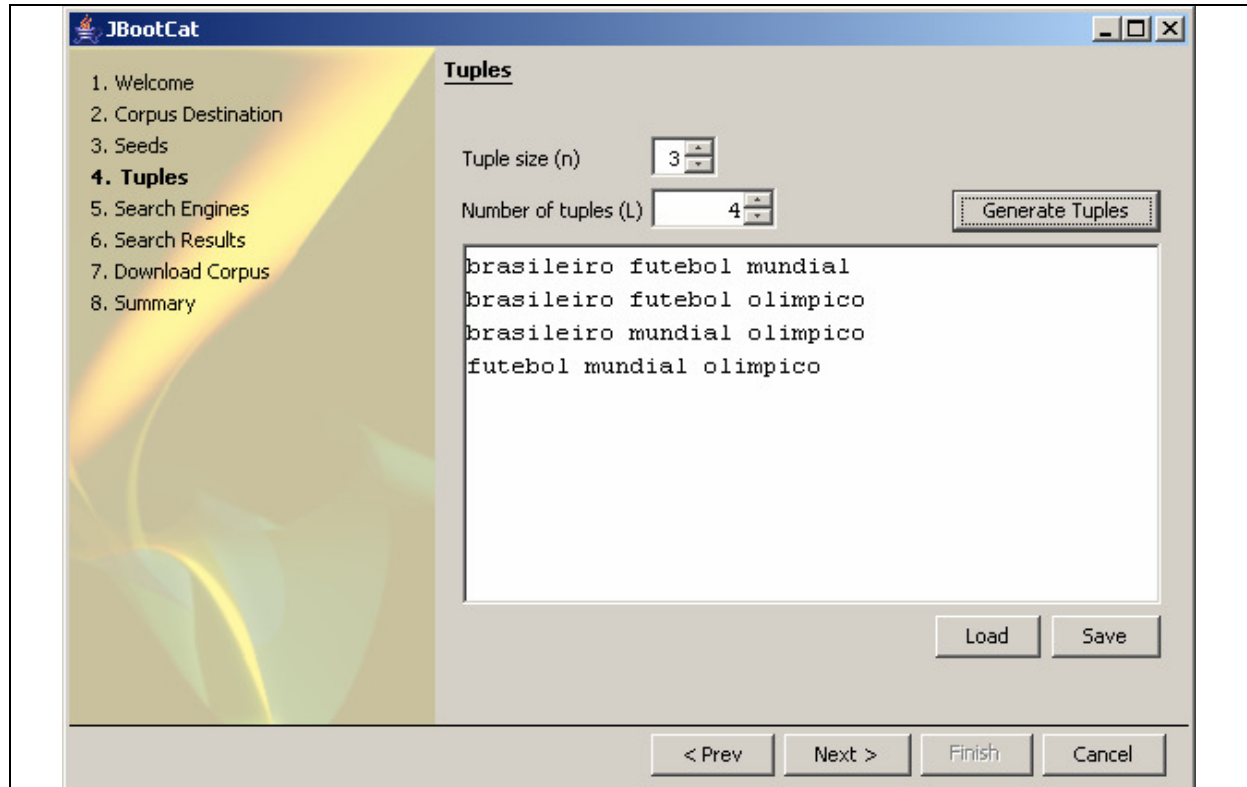
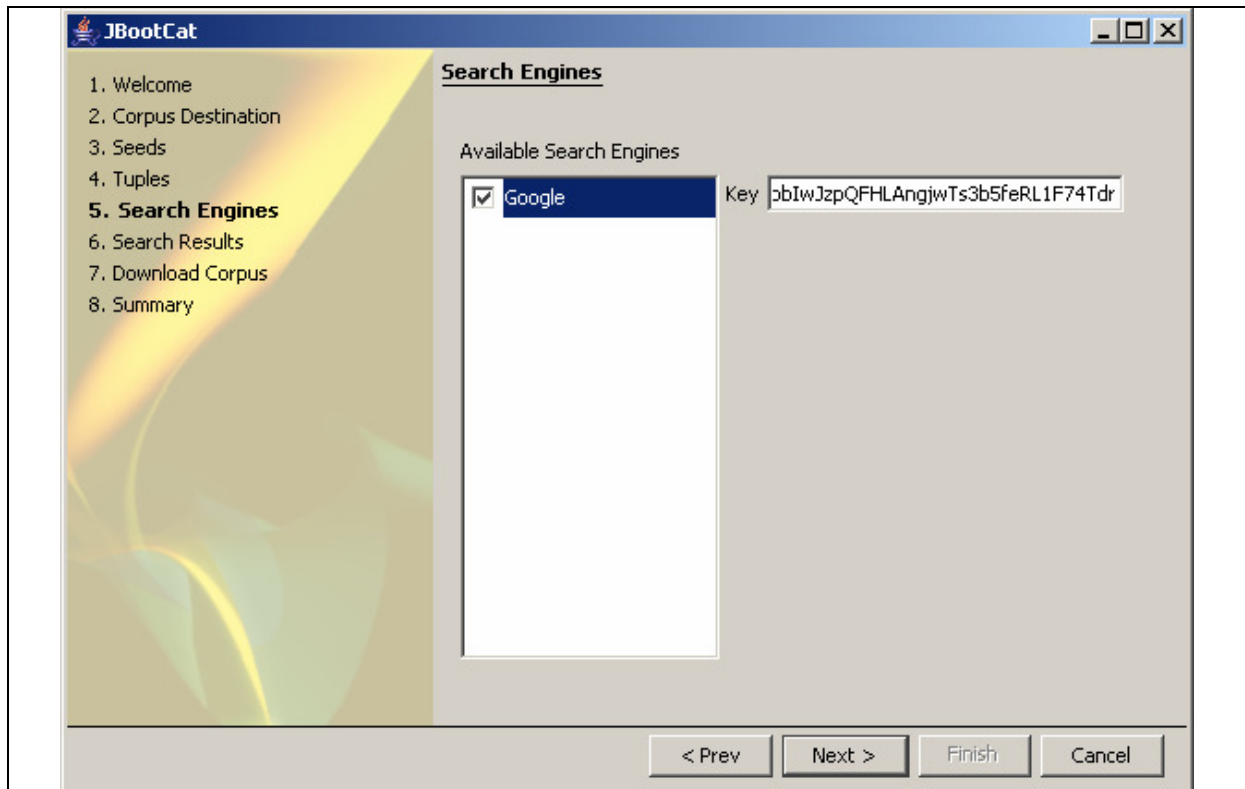
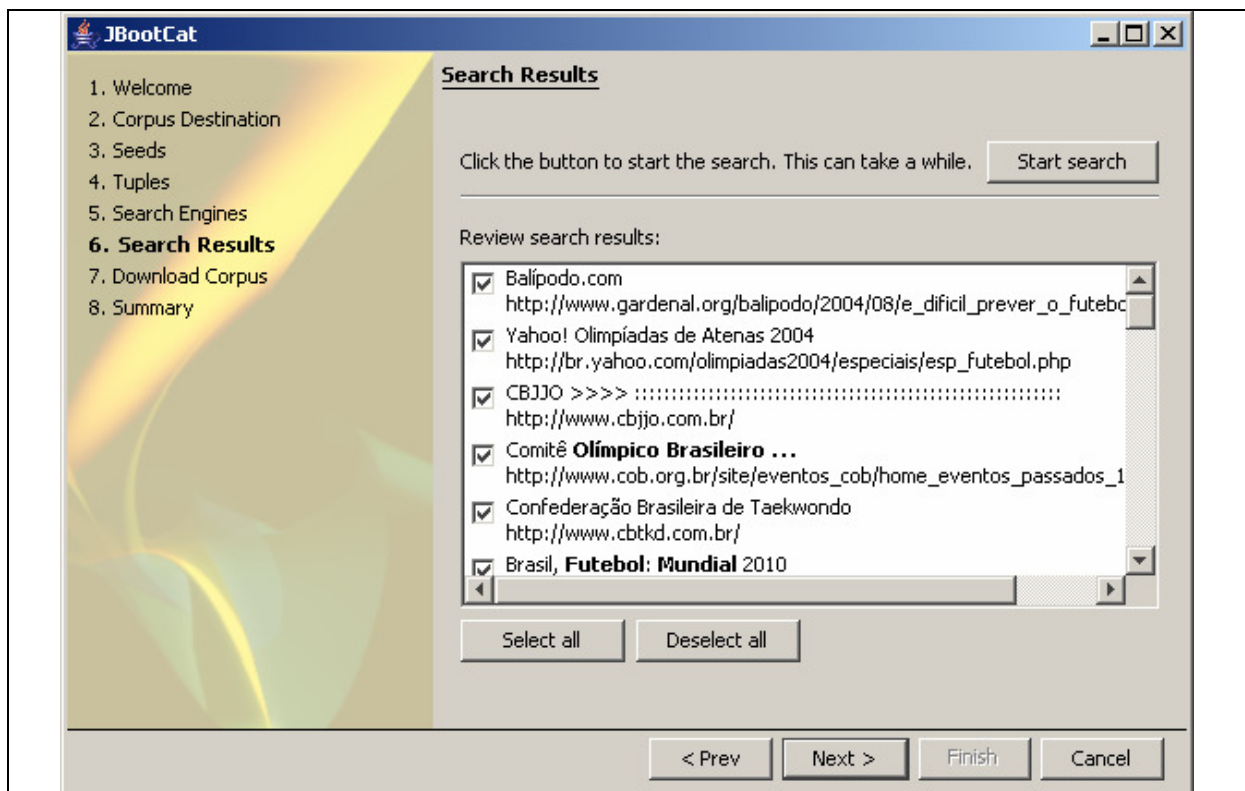


Figura 6.14 – Geração das tuplas para pesquisa





**Figura 6.15 – Escolha do motor de busca**



**Figura 6.16 – Resultado da pesquisa na web**

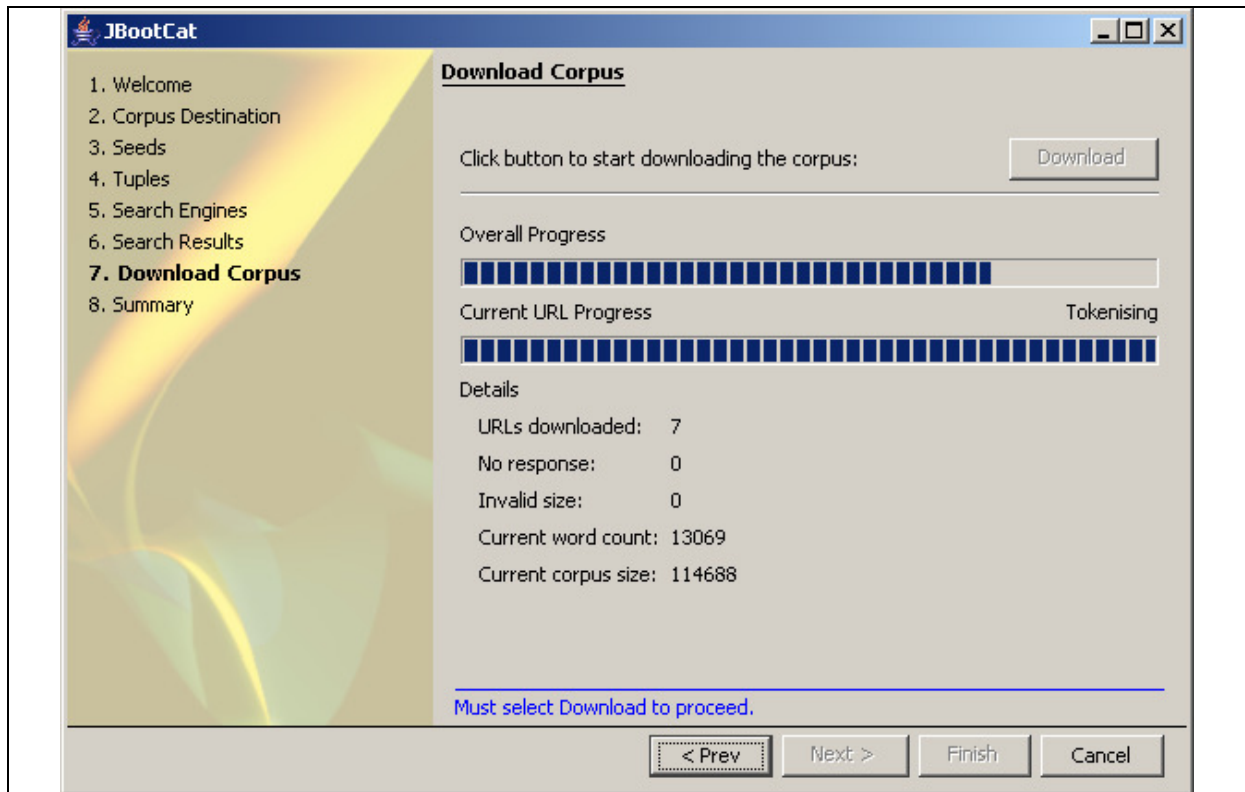


Figura 6.17 – Processamento da pesquisa pela ferramenta

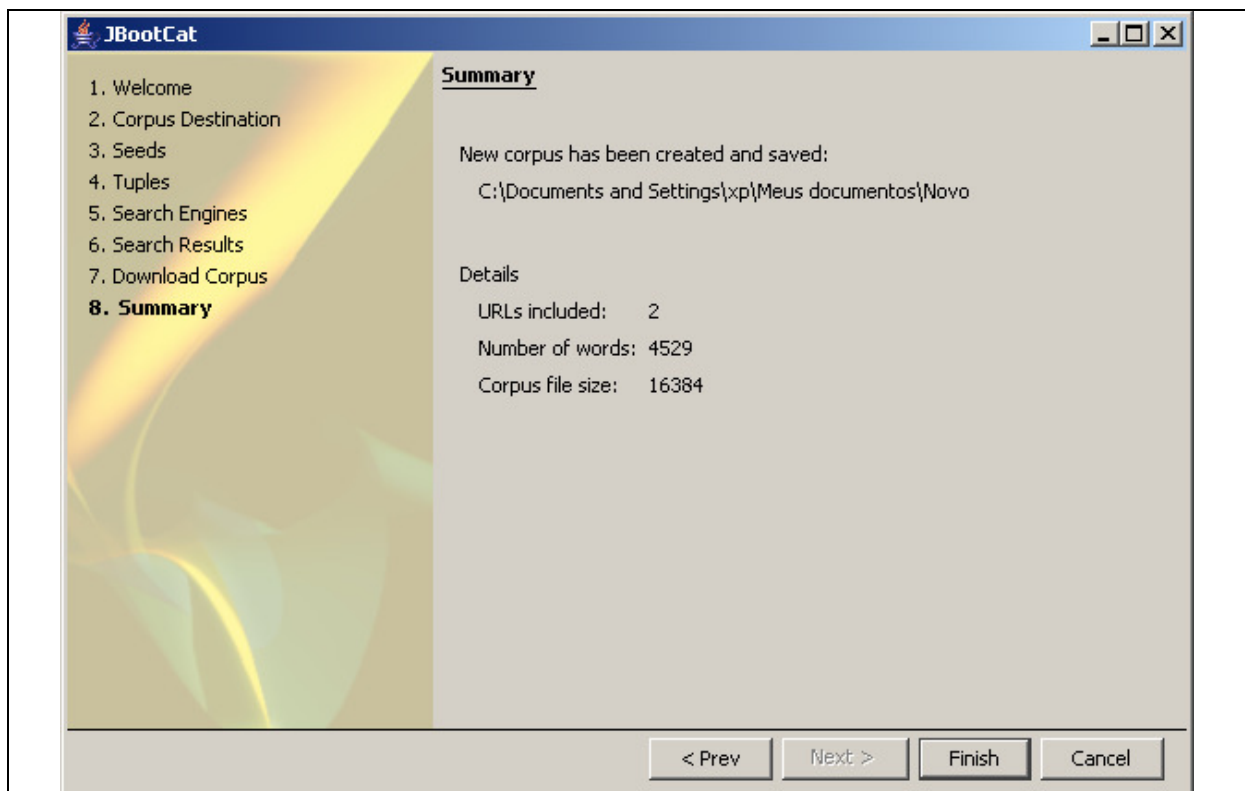


Figura 6.18 – Resumo das ações da ferramenta

### 6.5 API Apache Lucene

O *Apache Lucene* é um *software* de busca e ao mesmo tempo uma *API* de indexação de documentos, escrito na linguagem de programação Java. É um *software* de código aberto (*open-source*) da *Apache Software Foundation* licenciado através da licença *Apache* (LUCENE, 2007).

Para o *Lucene* não importa a origem dos dados, seu formato, ou mesmo a linguagem em que foi escrito, desde que esses dados possam ser convertidos para o formato de texto. Desta forma esta *API* pode ser utilizada para indexar e buscar dados gravados em diversas fontes de dados como, por exemplo, arquivos, páginas *web* em servidores remotos, documentos gravados no sistema de arquivos local, arquivos textos, documentos *Microsoft Word*, documentos *HTML* ou *PDF*, ou qualquer outro formato do qual possa ser extraído algum tipo de informação textual (LUCENE, 2007).

Nesta *API* podem ser encontrados algoritmos que executam algumas funcionalidades relacionadas com o PLN, tais como *stemming*, remoção de *stopwords*, contadores de ocorrências de palavras em textos, algoritmos que eliminam espaços brancos em textos, analisadores sintáticos, etc.

## 7 MODELAGEM DAS API'S DE MECANISMOS DE BUSCA

Este capítulo mostra a modelagem das *API's* de mecanismos de busca utilizados no decorrer do trabalho. A modelagem é apresentada através da utilização de casos de uso, para descrever de uma forma simples os passos que devem ser seguidos para utilizar determinada *API*, Diagramas de Classe, que mostram as principais classes envolvidas e por fim Diagramas de Sequência que exemplificam a utilização das classes.

### 7.1 Modelagem *Google API*

Nesta seção é mostrada a modelagem realizada para exemplificar a utilização da *Google API*. Primeiramente é mostrado o caso de uso da principal funcionalidade desta *API* para o trabalho, a pesquisa na *web*. Após é mostrado o Diagrama de Classes, mostrando as principais classes envolvidas neste processo e por fim o Diagrama de Sequência.

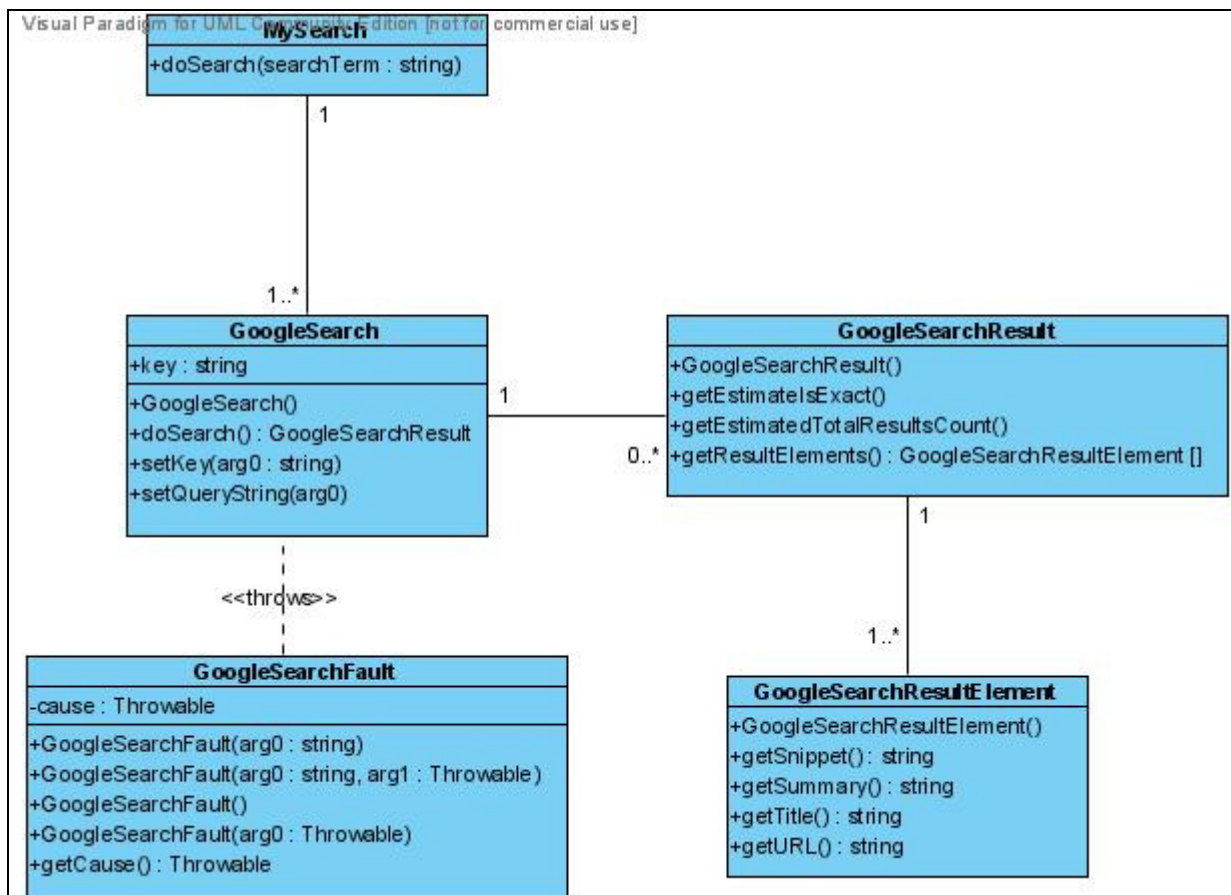
#### 7.1.1 Realizar consultas através da *Google API*

Neste caso de uso é mostrado o processo de pesquisa feito através da utilização da *API* do *Google*.

<b>CDU-001 – Realizar Consulta Através da <i>Google API</i></b>	
<b>Objetivo:</b>	Realizar consulta na <i>web</i> através da utilização da <i>API</i> do <i>Google</i>
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário (identificado como Ator) e <i>API</i> do <i>Google</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator informa sobre qual termo de pesquisa ele deseja realizar a pesquisa	
2. Ator informa a chave de pesquisa	
3. Sistema recebe o termo e a chave de pesquisa informada pelo Ator	
4. Sistema pesquisa o termo informado pelo Ator na base do <i>Google</i>	

5. Sistema exibe os resultados da pesquisa. Resultados exibidos: a) Título da Página b) URL da Página c) Resumo
6. Caso de uso se encerra
<b>Cenários Alternativos:</b>
Nenhum resultado encontrado:
<b>5.a.</b> Sistema não encontra nenhum resultado para o termo informado pelo Ator e informa ao mesmo que não foi possível encontrar resultados para a pesquisa na base de dados do <i>Google</i>

### 7.1.2 Diagrama de Classes da *Google API*



**Figura 7.1 – Diagrama de Classes da *Google API***

A classe “*MySearch*” é a responsável por instanciar as outras classes (*GoogleSearch*, *GoogleResultElement*, *GoogleSearchResultElement*, *GoogleSearchFault*) que serão utilizadas para realizar a consulta. Estas classes são chamadas através do método “*doSearch*”. Neste método é setada a chave de pesquisa, através do método “*setKey*” da classe “*GoogleSearch*”, e através do método “*setQueryString*”, da mesma classe, é setado o termo de pesquisa.

A classe “*GoogleSearch*” tem a finalidade de chamar os métodos que irão realizar a pesquisa na base do *Google*, através da *Google API*. O método “*doSearch*” desta classe é o responsável por realizar este processo.

A classe “*GoogleResultElement*” recebe a resposta vinda da pesquisa realizada na base de páginas do *Google*.

A classe “*GoogleSearchResultElement*” recebe todos os resultados da classe anterior e separa cada resultado em itens, que (Título da Página, URL, Resumo) são aqueles que o usuário visualizará como resultado da pesquisa. Os métodos que retornam estes itens são respectivamente “*getTitle*”, “*getUrl*” e “*getSummary*”. O método “*getSnippet*” retorna um pequeno texto sobre a página resultante da pesquisa.

A classe “*GoogleSearchFault*” é a responsável por mostrar algum eventual erro ocorrido durante o processo de pesquisa realizado na base do *Google*.

### 7.1.3 Diagrama de Seqüência da *Google API*

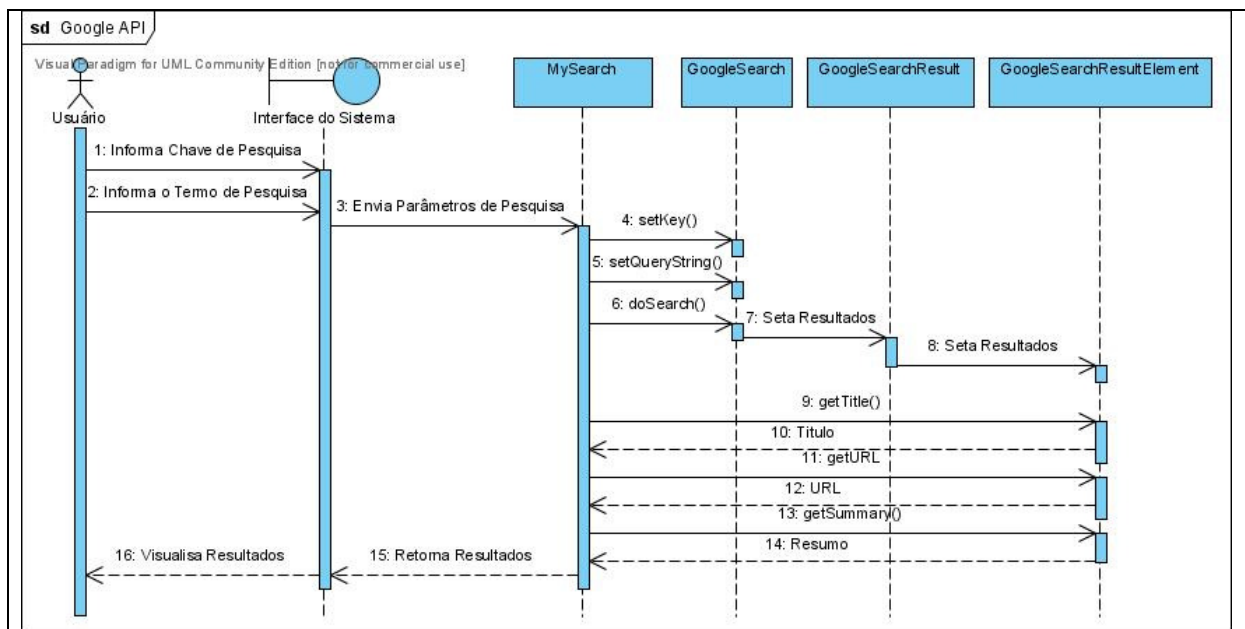


Figura 7.2 – Diagrama de Seqüência da *Google API*

## 7.2 Modelagem *Yahoo API*

Nesta seção, assim como na anterior, é mostrada a modelagem da *API* do *Yahoo* para realizar consultas na *web*.

### 7.2.1 Realizar consultas através da *Yahoo API*

Neste caso de uso é mostrado o processo de pesquisa feito através da utilização da *API* do *Yahoo*.

<b>CDU-001 – Realizar Consulta Através da <i>Yahoo API</i></b>	
<b>Objetivo:</b>	Realizar consulta na <i>web</i> através da utilização da <i>API</i> do <i>Yahoo</i>
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário (identificado como Ator) e <i>API</i> do <i>Yahoo</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator informa sobre qual termo de pesquisa ele deseja realizar a pesquisa	
2. Ator informa a chave de pesquisa	
3. Sistema recebe o termo e a chave de pesquisa informada pelo Ator	
4. Sistema pesquisa o termo informado pelo Ator na base do <i>Yahoo</i>	
5. Sistema exhibe os resultados da pesquisa. Resultados exibidos: a) Titulo da Página b) URL da Página c) Resumo	
6. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	
Nenhum resultado encontrado:	
5.a. Sistema não encontra nenhum resultado para o termo informado pelo Ator e informa ao mesmo que não foi possível encontrar resultados para a pesquisa na base de dados do <i>Yahoo</i>	

## 7.2.2 Diagrama de Classes da *Yahoo API*

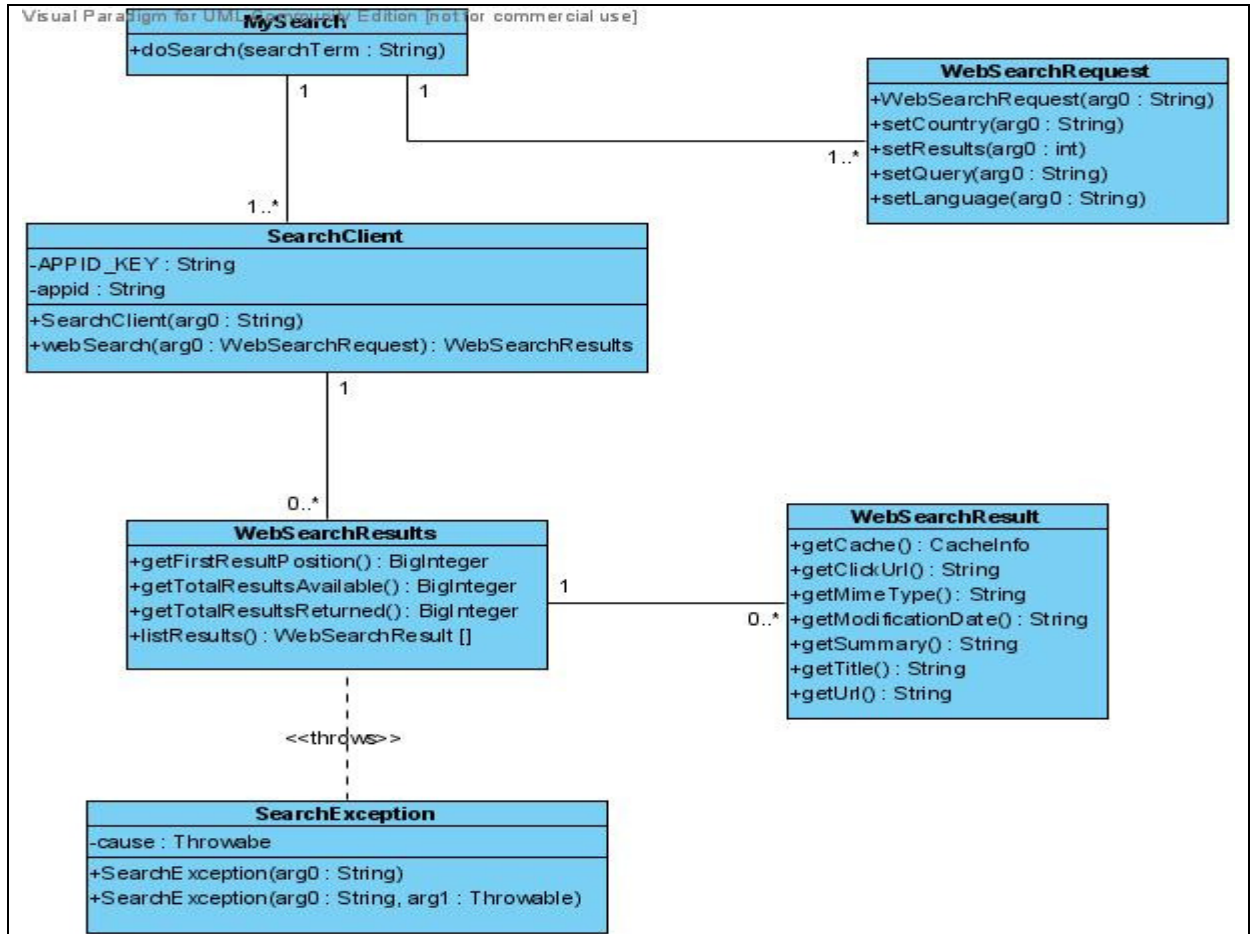


Figura 7.3 – Diagrama de Classes da *Yahoo API*

A classe “*MySearch*” através do seu método único, chamado “*doSearch*”, é a responsável por iniciar o processo de pesquisa, pois através dela as outras classes são instanciadas.

A primeira classe a ser instanciada é a “*SearchClient*”, considerada como sendo a classe principal para realizar pesquisas (na base do *Yahoo*) através da *API do Yahoo*. O método “*webSearch*” desta classe é o responsável por realizar a consulta na base.

A classe “*WebSearchRequest*” é a classe que recebe o termo de pesquisa do usuário, que é utilizado no método “*webSearch*” da classe “*SearchClient*” para realizar a consulta.

A classe “*WebSearchResults*” é a responsável por receber o retorno das consultas realizadas, enquanto a classe “*WebSearchResult*” recebe todos os resultados da classe anterior e separa cada resultado em itens, estes itens (Título da Página, URL, Resumo) são os que o



usuário visualizará como resultado final da pesquisa. Os métodos que retornam estes itens são respectivamente “*getTitle*”, “*getUrl*” e “*getSummary*”.

A classe “*SearchException*” é a responsável por mostrar algum eventual erro ocorrido durante o processo de pesquisa. Estes erros podem ocorrer quando, por exemplo, nenhuma pagina é retornada pela *API*.

### 7.2.3 Diagrama de Seqüência da *Yahoo API*

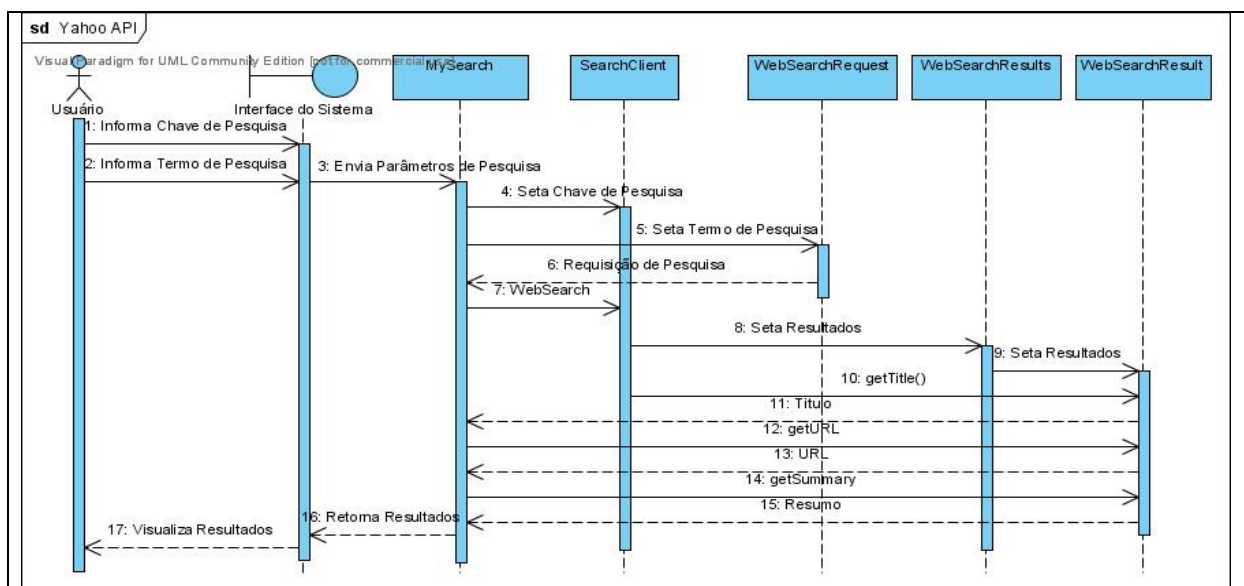


Figura 7.4 – Diagrama de Seqüência da *Yahoo API*

### 7.3 Modelagem da *Technorati API*

Nesta seção é mostrada uma breve modelagem da *API* do *Technorati*, onde é destacada a principal funcionalidade realizada por esta *API*, a pesquisa na base de *blogs* do *Technorati*.

#### 7.3.1 Realizar consultas através da *Technorati API*

Neste caso de uso é demonstrado o processo de pesquisa na base do *Technorati* para encontrar páginas que contenham o termo informado pelo usuário.

#### CDU-001 – Realizar Consulta Através da *Technorati API*

<b>Objetivo:</b>	Realizar consulta na <i>web</i> através da utilização da <i>API</i> do <i>Technorati</i>
------------------	--

<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário (identificado como Ator) e API do <i>Technorati</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator informa sobre qual termo de pesquisa ele deseja realizar a pesquisa	
2. Ator informa a chave de pesquisa	
3. Sistema recebe o termo e a chave de pesquisa informada pelo Ator	
4. Sistema pesquisa o termo informado pelo Ator na base de <i>blogs</i> do <i>Technorati</i>	
5. Sistema exhibe os resultados da pesquisa. Resultados exibidos: a) Título da Página b) URL da Página	
6. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	
Nenhum resultado encontrado:	
5.a. Sistema não encontra nenhum resultado para o termo informado pelo Ator e informa o mesmo que não foi possível encontrar resultados para a pesquisa na base de dados	

### 7.3.2 Diagrama de classes da *Technorati* API

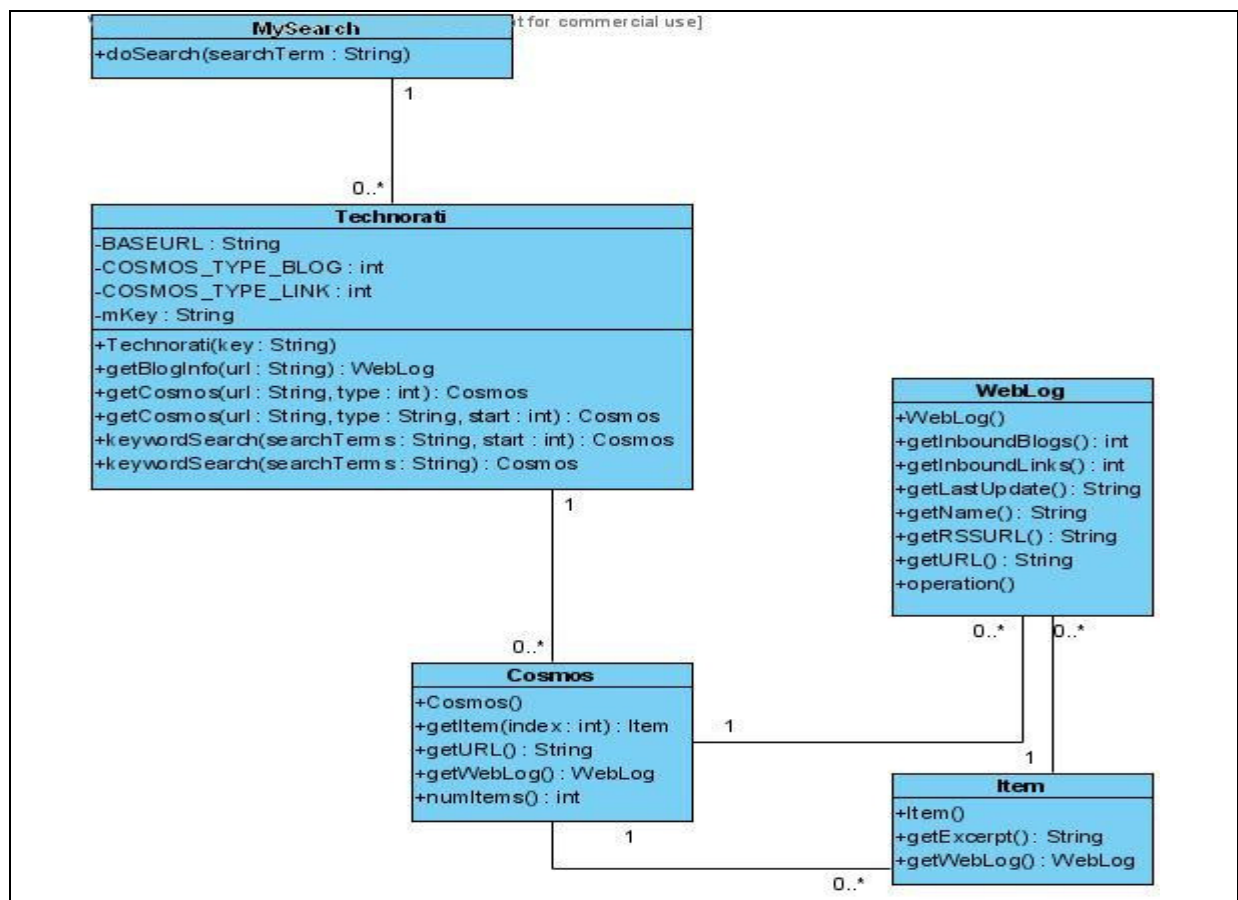


Figura 7.5 – Diagrama de Classes da *Technorati* API

A classe “*MySearch*”, assim como nas outras *API*’s modeladas, é a classe responsável por instanciar as demais classes do diagrama. O método “*doSearch*” é quem inicia o processo de pesquisa.

A classe “*Technorati*” é quem realiza o processo de pesquisa na base de *blogs* do *Technorati* através do método “*keyWordSearch*”. Este método é quem recebe o termo a ser pesquisado.

A classe “*Cosmos*” é a classe que recebe o retorno da pesquisa realizada pelo usuário na base do *Technorati*, através do método anteriormente citado.

Através do método “*getItem*” da classe “*Cosmos*” a classe “*Item*” recebe os itens retornados pela pesquisa.

A classe “*WebLog*”, através do método “*getWebLog*”, da classe anterior recebe os resultados que serão mostrados para o usuário, tais como Título da Página (método “*getName*”) e URL (método “*getURL*”). O Resumo é retornado pela classe “*Item*” (método “*getExcerpt*”).

### 7.3.3 Diagrama de Seqüência da *Technorati API*

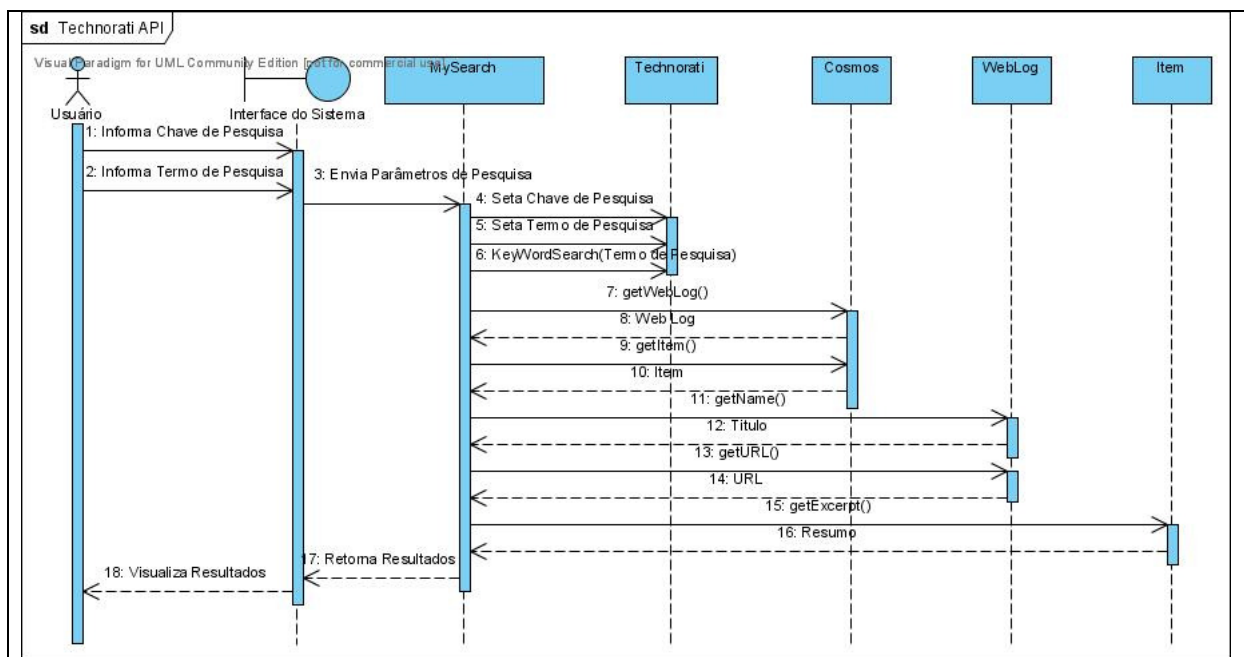


Figura 7.6 – Diagrama de Seqüência da *Technorati API*

## 8 MODELAGEM DAS FERRAMENTAS DE PLN

Este capítulo tem como objetivo mostrar a modelagem de casos de uso das ferramentas de PLN estudadas no decorrer deste trabalho. E por fim a modelagem de algumas funcionalidades que podem ser implementadas através da utilização da *API Lucene*.

### 8.1 Modelagem de Casos de Uso da Ferramenta *WebJspell*

Nesta seção são mostrados os casos de uso extraídos a partir da utilização da ferramenta *WebJspell*.

#### 8.1.1 Informar a frase que será processada e alguns parâmetros iniciais

Este é o caso de uso inicial, onde será informada a frase para posterior análise morfológica e os parâmetros iniciais solicitados pelo sistema.

<b>CDU-001 – Informar a frase que será processada e alguns parâmetros iniciais</b>	
<b>Objetivo:</b>	Informar a frase que será analisada morfológicamente pelo analisador morfológico
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário do <i>WebJspell</i> (identificado como Ator) e <i>WebJspell</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator acessa o Sistema através do endereço “ <a href="http://linguateca.di.uminho.pt/webjspell/jsol.pl">http://linguateca.di.uminho.pt/webjspell/jsol.pl</a> ”	
2. Ator a frase que deseja que o Sistema faça a análise morfológica	
3. Ator seleciona o dicionário que pretende utilizar. Existem três opções: a) Português b) Latim (Incompleto); c) Inglês.	
4. Ator seleciona sobre qual nível de detalhe deseja que o Sistema retorne os resultados da análise. Existem duas opções: a) Normal	

b) Resumido
<b>5.</b> Ator solicita que o Sistema realize a análise morfológica através do botão “Analisar”
<b>6.</b> Caso de uso se encerra
<b>Cenários Alternativos:</b>

### 8.1.2 Análise morfológica da frase informada

Neste caso de uso o sistema realizará a análise morfológica da frase informada pelo usuário, retornando a análise de cada palavra separadamente.

<b>CDU-002 – Análise morfológica da frase informada</b>	
<b>Objetivo:</b>	Realizar a análise morfológica da frase informada pelo usuário
<b>Pré-Condição:</b>	CDU-001
<b>Atores:</b>	Usuário do <i>WebJspell</i> (identificado como Ator) e <i>WebJspell</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
<b>1.</b> Sistema realiza a análise morfológica de cada palavra pertencente a frase informada pelo Ator	
<b>2.</b> Ator a frase que deseja que o Sistema faça a análise morfológica	
<b>3.</b> Sistema retorna a análise de cada palavra mostrando suas particularidades. Particularidades como: a) Categoria; b) Número; c) Caso latino; d) Pessoa; e) Lema; f) Tempo (verbo); g) Transitividade (verbo); h) Gênero.	
<b>4.</b> Sistema retorna algumas opções para o usuário. As opções retornadas são as seguintes: a) Sugerir análise; b) Remover do Dicionário; c) Flexionar; d) Exemplos das frases; e) Ocorrências; f) Conjugação verbal (verbos)	
<b>5.</b> Ator solicita que o Sistema realize a análise morfológica através do botão “Analisar”	
<b>6.</b> Caso de uso se encerra	
<b>Cenários Alternativos:</b>	
Palavra não encontrada no dicionário:	

**2.a.** Sistema não encontra a palavra no dicionário e não mostra nenhuma particularidade da mesma. Sistema retorna algumas opções para o usuário: a) Adicionar palavra ao dicionário; b) Corrigir palavra; c) Ocorrências.

### 8.1.3 Sugerir Análise

Nesta opção o usuário pode enviar para a equipe de desenvolvimento do *WebJspell* uma sugestão de análise morfológica diferente da exibida.

<b>CDU-003 – Sugerir Análise</b>	
<b>Objetivo:</b>	Sugerir uma análise morfológica diferente para a palavra analisada
<b>Pré-Condição:</b>	CDU-002
<b>Atores:</b>	Usuário do <i>WebJspell</i> (identificado como Ator) e <i>WebJspell</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator redige a nova análise morfológica para a palavra	
2. Ator opcionalmente informa seus dados: a) Nome; b) E-mail.	
3. Ator submete sua análise para o Sistema através do botão “Submeter”	
4. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	

### 8.1.4 Remover do dicionário

Nesta opção o usuário pode enviar para a equipe de desenvolvimento do *WebJspell* uma sugestão de remoção da palavra selecionada.

<b>CDU-004 – Remover do dicionário</b>	
<b>Objetivo:</b>	Remover análise morfológica da palavra selecionada do sistema
<b>Pré-Condição:</b>	CDU-002
<b>Atores:</b>	Usuário do <i>WebJspell</i> (identificado como Ator) e <i>WebJspell</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator sugere a remoção da palavra	

<b>2.</b> Ator opcionalmente informa seus dados: a) Nome; b) E-mail.
<b>3.</b> Ator seleciona o dicionário ao qual quer excluir a palavra. Existem três opções: a) Português; b) Latim (Incompleto); c) Inglês.
<b>4.</b> Ator submete sua sugestão de exclusão para o sistema
<b>5.</b> Caso de uso se encerra
<b>Cenários Alternativos:</b>

### 8.1.5 Flexionar a palavra

Nesta opção o usuário solicita que o sistema mostre todas as flexões disponíveis para a palavra selecionada.

<b>CDU-005 – Flexionar a palavra</b>	
<b>Objetivo:</b>	Mostrar ao usuário as formas flexionadas disponíveis para a palavra selecionada
<b>Pré-Condição:</b>	CDU-002
<b>Atores:</b>	Usuário do <i>WebJspell</i> (identificado como Ator) e <i>WebJspell</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
<b>1.</b> Ator solicita ao Sistema que ele retorne todas as formas flexionadas disponíveis para a palavra selecionada	
<b>2.</b> Sistema retorna em um quadro contendo as formas flexionadas para a palavra selecionada	
<b>3.</b> Caso de uso se encerra	
<b>Cenários Alternativos:</b>	

### 8.1.6 Exemplos de frases

Nesta opção o usuário solicita ao sistema que ele busque em sua base alguns exemplos de frases onde a palavra selecionada está inserida.

<b>CDU-006 – Exemplos de frases</b>	
<b>Objetivo:</b>	Mostrar ao usuário alguns exemplos de frases contendo a palavra selecionada

<b>Pré-Condição:</b>	CDU-002
<b>Atores:</b>	Usuário do <i>WebJspell</i> (identificado como Ator) e <i>WebJspell</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator solicita ao sistema que ele retorne alguns exemplos de frase onde a palavra selecionada esta presente	
2. Sistema busca em sua base exemplos de frases	
3. Sistema exibe um quadro contendo as frases encontradas destacando a palavra selecionada	
4. Sistema exibe um resumo com as ações executadas. Estes resumos contem as seguintes informações: a) Número de exemplos retornados; b) Quantidade de frases presentes na base do Sistema que contem a palavra selecionada	
5. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	
Exemplos não encontrados:	
3.a. Sistema não encontra frases na sua base onde a palavra selecionada esta presente e exibe a seguinte mensagem: “ <b>Não foi possível obter exemplos de frases</b> ”	

### 8.1.7 Ocorrências da palavra

Nesta opção o usuário solicita ao sistema que ele exiba a quantidade de ocorrências da palavra na base.

<b>CDU-007 – Ocorrências da palavra</b>	
<b>Objetivo:</b>	Mostrar ao usuário a quantidade de ocorrências da palavra selecionada na base do sistema
<b>Pré-Condição:</b>	CDU-002
<b>Atores:</b>	Usuário do <i>WebJspell</i> (identificado como Ator) e <i>WebJspell</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator solicita ao sistema que ele retorne a quantidade de ocorrências da palavra selecionada na base do Sistema	
2. Sistema busca na sua base a quantidade de ocorrências para a palavra selecionada	
3. Sistema exibe um quadro contendo a quantidade de ocorrências encontradas para a palavra selecionada	
4. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	



Exemplos não encontrados:
<b>3.a.</b> Sistema não encontra nenhuma ocorrência para palavra selecionada e exibe a seguinte mensagem: “ <b>Não foi encontrada qualquer ocorrência</b> ”

### 8.1.8 Conjugação verbal

Nesta opção o usuário solicita ao sistema que ele exiba a conjugação verbal.

<b>CDU-008 – Conjugação verbal</b>	
<b>Objetivo:</b>	Mostrar ao usuário a conjugação verbal para a palavra
<b>Pré-Condição:</b>	CDU-002
<b>Atores:</b>	Usuário do <i>WebJspell</i> (identificado como Ator) e <i>WebJspell</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator solicita ao Sistema que ele exiba a conjugação verbal da palavra	
2. Sistema consulta a base para buscar a conjugação verbal	
3. Sistema exibe um quadro contendo todas as conjugações verbais para a palavra	
4. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	

## 8.2 Modelagem de Casos de Uso da Ferramenta *Snowball*

Nesta seção serão mostrados os casos de uso extraídos através da utilização da ferramenta *Snowball*.

### 8.2.1 Informar as palavras para realizar *stemming* ou lematização

Este é o caso de uso inicial, onde serão informadas as palavras para posterior *stemming* ou lematização.

<b>CDU-001 – Informar as palavras para realizar <i>stemming</i> ou lematização</b>	
<b>Objetivo:</b>	Informar as palavras para posterior <i>stemming</i> ou lematização
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário do <i>Snowball</i> (identificado como Ator) e <i>Snowball</i> (identificado como Sistema)

<b>Cenário Principal:</b>
1. Ator acessa o sistema através do endereço “http://www.snowball.tartarus.org/demo.php”
2. Ator informa as palavras que deseja que o sistema faça o <i>stemming</i> ou lematização
3. Ator solicita que o sistema realize o <i>stemming</i> ou lematização através do botão “Submit”
4. Caso de uso se encerra
<b>Cenários Alternativos:</b>

### 8.2.2 Realizar *stemming* ou lematização

Neste caso de uso o sistema realizará o *stemming* ou lematização das palavras informadas pelo usuário.

<b>CDU-002 – Realizar <i>stemming</i> ou lematização</b>	
<b>Objetivo:</b>	Realizar o <i>stemming</i> ou lematização das palavras informadas pelo usuário
<b>Pré-Condição:</b>	CDU-001
<b>Atores:</b>	<i>Snowball</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Sistema realiza o processo de <i>stemming</i> ou lematização das palavras informadas pelo usuário	
2. Sistema exibe o resultado do processo mostrando as palavras informadas pelo usuário, seguido do resultado do <i>stemming</i> ou lematização	
3. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	

### 8.3 Modelagem de Casos de Uso da Ferramenta *Flip-Online*

Nesta seção serão mostrados os casos de uso extraídos através da utilização do *Flip-Online*.

#### 8.3.1 Informar a frase que será processada e os parâmetros iniciais

Este é o caso de uso inicial, onde será informada a frase para posterior análise sintática e os parâmetros iniciais solicitados pelo sistema.

<b>CDU-001 – Informar a frase que será processada e os parâmetros iniciais</b>	
<b>Objetivo:</b>	Informar a frase que será analisada morfológicamente pelo analisador morfológico
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário do <i>Flip-Online</i> (identificado como Ator) e <i>Flip-Online</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator acessa o sistema através do endereço “ <a href="http://www.flip.pt/FLiPOnline/tabid/96/Default.aspx">http://www.flip.pt/FLiPOnline/tabid/96/Default.aspx</a> ”	
2. Ator informa a frase que deseja que o sistema faça a análise sintática	
3. Ator seleciona sobre qual língua o sistema irá processar a frase. As opções são: a) Português (Portugal) b) Português (Brasil)	
4. Ator solicita que o sistema realize a análise sintática através do botão “Verificar Texto”	
5. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	

### 8.3.2 Análise sintática da frase informada

Neste caso de uso o sistema realizará a análise sintática da frase informada pelo usuário.

<b>CDU-002 – Análise sintática da frase informada</b>	
<b>Objetivo:</b>	Realizar a análise sintática da frase informada pelo usuário
<b>Pré-Condição:</b>	CDU-001
<b>Atores:</b>	<i>Flip-Online</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Sistema realiza análise sintática da frase informada pelo Ator	
2. Sistema sublinha as palavras incorretas sintaticamente	
3. Sistema exibe a quantidade de erros encontrados na frase	
4. Sistema exibe opções de correção de cada palavra sintaticamente incorreta. As opções são as seguintes: a) Exibe a frase “ <b>Verifique se quer manter</b> ” a palavra sintaticamente incorreta; b) Ignorar; c) Alterar; d) Cancelar.	

<b>5. Caso de uso se encerra</b>
<b>Cenários Alternativos:</b>

#### 8.4 Modelagem de Casos de Uso da Ferramenta *JBootCat*

Nesta seção serão mostrados os casos de uso extraídos através da utilização da ferramenta *JBootCat*.

##### 8.4.1 Informar arquivo para armazenamento do *Corpus*

Antes de realizar qualquer pesquisa o usuário deverá informar ao sistema em qual arquivo e diretório o *corpus* que será gerado deverá ser armazenado.

<b>CDU-001 – Informar Arquivo para Armazenamento de <i>Corpus</i></b>	
<b>Objetivo:</b>	Informar arquivo e diretório onde será armazenado o <i>corpus</i> resultante da pesquisa
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário do <i>JBootCat</i> (identificado como Ator) e <i>JBootCat</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator seleciona o nome do arquivo e o diretório onde será gravado o <i>corpus</i> resultante da pesquisa que será realizada	
2. Ator confirma as informações para o Sistema	
3. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	

##### 8.4.2 Informar termos de pesquisa

Este é considerado como sendo o caso de uso inicial do sistema, pois é através dele que são informados quais os termos que o usuário deseja pesquisar para construir seu *corpus*. Este caso de uso é extremamente curto, porém de muita importância no decorrer da modelagem.

<b>CDU-002 – Informar termos de pesquisa</b>

<b>Objetivo:</b>	Informar quais termos de pesquisa serão utilizados na pesquisa
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário do <i>JBootCat</i> (identificado como Ator) e <i>JBootCat</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator informa quais termos ele deseja pesquisar através do Sistema. Os termos são palavras a serem pesquisadas na <i>web</i> através das <i>API's</i> dos mecanismos de busca	
2. Ator confirma os termos informados	
3. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	

### 8.4.3 Montar tuplas de pesquisa

Neste caso de uso o usuário poderá escolher quantas tuplas<sup>20</sup> de pesquisa ele deseja gerar com os termos por ele informados. Gerar tuplas significa dividir os termos em listas distintas, onde cada linha desta lista conterà uma combinação dos termos pesquisados.

<b>CDU-003 – Montar tuplas de pesquisa</b>	
<b>Objetivo:</b>	Montar tuplas de pesquisa com os termos informados pelo usuário
<b>Pré-Condição:</b>	CDU-002
<b>Atores:</b>	Usuário do <i>JBootCat</i> (identificado como Ator) e <i>JBootCat</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator informa quantas tuplas ele deseja gerar com os termos por ele informados	
2. Ator solicita a geração das tuplas para o Sistema	
3. Sistema processa os termos informados pelo ator para montar as tuplas de pesquisa	
4. Sistema exibe as tuplas de pesquisa geradas para o Ator	
5. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	
Número de tuplas maior que o de termos:	
2.a. O Sistema informa ao Ator que o número de tuplas de pesquisa solicitadas para geração excede o número de termos informados pelo Ator, e solicita a redução da quantidade de tuplas de pesquisa	

<sup>20</sup> Tupla - A tupla é uma seqüência, como a lista: armazena um conjunto de elementos acessíveis por um índice inteiro. A tupla é imutável; uma vez criada, não pode ser modificada.

Muitas tuplas de pesquisa solicitadas:
<b>2.b.</b> O Sistema informa ao Ator que foram solicitadas muitas tuplas de pesquisa para o número de combinações possíveis com a quantidade de termos informados e solicita que o Ator informe um novo número de tuplas de pesquisa

#### 8.4.4 Realizar pesquisa na *web*

Neste ponto o usuário solicitará que a pesquisa por ele informada (com base nos termos de pesquisa e tuplas de pesquisa geradas) seja realizada. O sistema por sua vez consultará a *web* através dos mecanismos de busca e trará ou não os resultados encontrados.

<b>CDU-004 – Realizar pesquisa na <i>web</i></b>	
<b>Objetivo:</b>	Realizar consultas de termos informados pelo usuário na <i>web</i> , através da utilização de <i>API's</i> de Mecanismos de Busca
<b>Pré-Condição:</b>	CDU-003
<b>Atores:</b>	Usuário do <i>JBootCat</i> (identificado como Ator) e <i>JBootCat</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Com base nas tuplas de pesquisa geradas o Ator executa a pesquisa	
2. Ator realiza a pesquisa na <i>web</i>	
3. Sistema verifica se o link encontrado se refere a uma página contendo <i>HTML</i>	
4. Sistema retorna os resultados para o usuário contendo: a) O título da página; b) A <i>URL</i> da página; c) Um pequeno resumo da página.	
5. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	
Pesquisa não retorna nenhum resultado:	
2.a. O Sistema informa o Ator que não foram encontrados resultados para os termos pesquisados	
Página encontrada não é <i>HTML</i> :	
3.a. O Sistema encontra nos resultados retornados pela pesquisa um link que não aponta para uma página <i>HTML</i> (pode ser um link para um arquivo PDF <sup>21</sup> ). O Sistema verifica o próximo link encontrado	

#### 8.4.5 Selecionar resultados da pesquisa

<sup>21</sup> PDF – Do Inglês Portable Document Format (Formato de Documento Portátil). É um formato de arquivo desenvolvido pela Adobe para representar documentos de maneira independente do aplicativo, hardware, e sistema operacional usados para criá-los. Um PDF pode descrever documentos que contenham texto, gráficos e imagens num formato independente do dispositivo e resolução. Disponível em: <http://www.adobe.com/br/products/acrobat/adobepdf.html>.

Após ter informado os termos, gerado as tuplas de pesquisa e solicitado a consulta o usuário deverá escolher com quais itens resultantes ele deseja construir o *corpus*. Com base nestes itens selecionados o sistema irá, num passo posterior, gerar o *corpus*.

<b>CDU-005 – Selecionar resultados da pesquisa</b>	
<b>Objetivo:</b>	Selecionar quais itens resultantes da pesquisa de termos devem estar presentes na construção do <i>corpus</i>
<b>Pré-Condição:</b>	CDU-004
<b>Atores:</b>	Usuário do <i>JBootCat</i> (identificado como Ator) e <i>JBootCat</i> (identificado como Sistema)
<b>Cenário Principal:</b>	
1. Ator seleciona com quais itens resultantes da pesquisa de termos ele deseja montar o <i>corpus</i>	
2. Sistema processa os itens escolhidos	
3. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	

#### 8.4.6 Gerar corpus

Após ter completado todos os casos de uso anteriores o usuário poderá solicitar ao sistema a construção do *corpus* baseado nos itens por ele escolhidos. Quando o usuário solicita a geração do *corpus* o sistema acessa cada página selecionada e gera os resumos contendo caracteres desnecessários para o *corpus* gerado, como por exemplo, marcações de linguagens *web*, como *HTML* e *JavaScript*. Não sendo estas marcações importantes para o *corpus* do usuário elas devem ser retiradas do arquivo gerado. Por fim o *corpus* gerado será armazenado em um arquivo previamente informado pelo usuário, no primeiro caso de uso.

<b>CDU-006 – Gerar <i>corpus</i></b>	
<b>Objetivo:</b>	Gerar o <i>corpus</i> com base nos itens escolhidos pelo usuário
<b>Pré-Condição:</b>	CDU-005
<b>Atores:</b>	<i>JBootCat</i> (identificado como Sistema)
<b>Cenário Principal:</b>	

1. Com base nos itens selecionados o Sistema acessa as páginas através das URL's contidas nos itens
2. Sistema gera resumos de cada página
3. Sistema realiza a limpeza do resumo (retira marcações de linguagens de programação)
4. Sistema armazena o <i>corpus</i> no arquivo previamente escolhido pelo Ator (CDU-001)
5. Sistema apresenta um resumo com detalhes das ações realizadas contendo: a) Contagem do número de palavras que estão presentes no <i>corpus</i> ; b) Erros encontrados; c) Quantidade de <i>URL's</i> baixadas com sucesso; d) Tamanho do <i>corpus</i> gerado.
6. Caso de uso se encerra
<b>Cenários Alternativos:</b>

## 8.5 Modelagem da *API Lucene*

Nesta seção é mostrada a modelagem da *API Lucene*, onde são destacadas algumas funcionalidades contidas na mesma.

### 8.5.1 Algoritmo de *Stemming* da *API Lucene*

Nesta seção é apresentado o algoritmo de *Stemming* que pertence a *API Lucene*. Através dele o usuário pode realizar o processo de radicalização de palavras. A seguir é apresentada a modelagem realizada para demonstrar sua utilização.

#### 8.5.1.1 Realizar *Stemming* em um texto

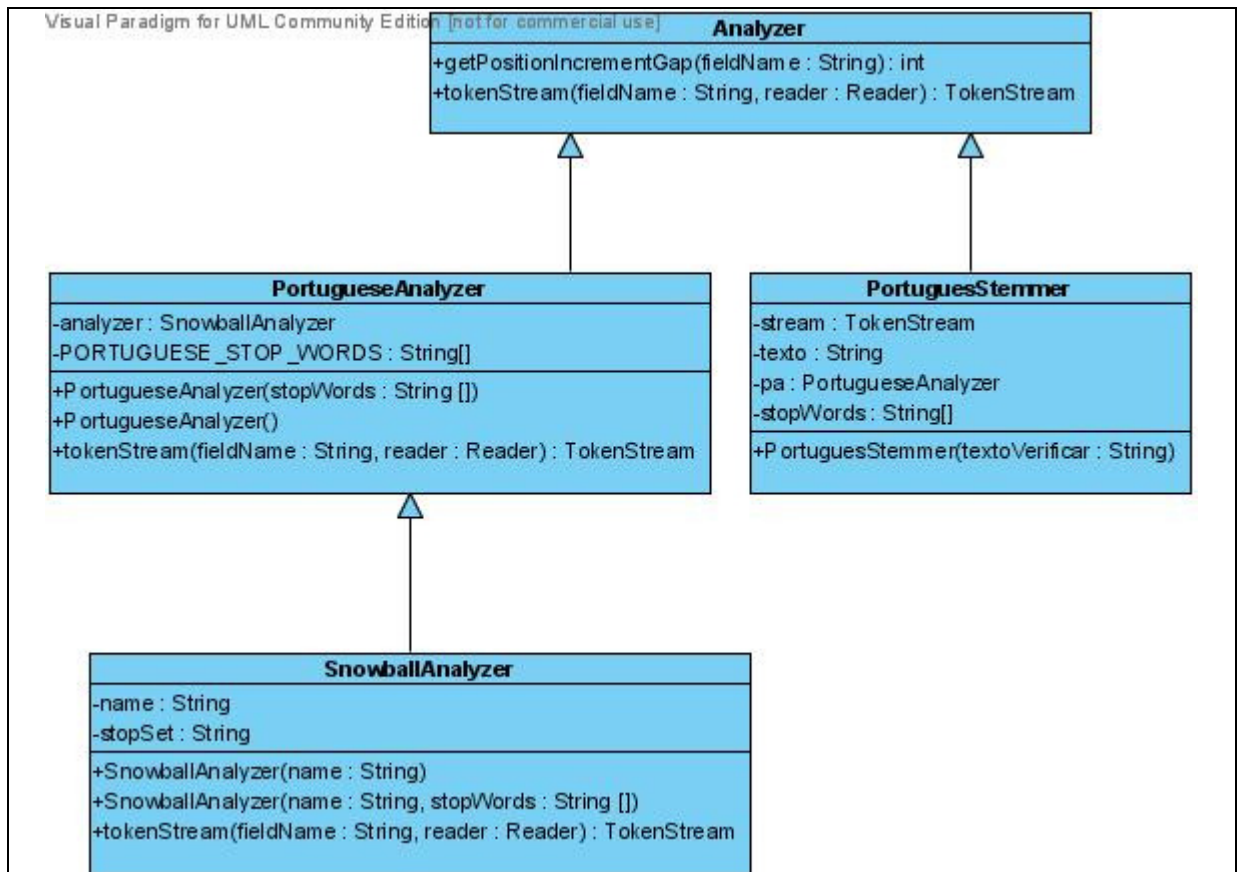
Nesta subseção é mostrado o principal caso de uso extraído através do desenvolvimento de um exemplo utilizando o algoritmo de *Stemming*.

<b>CDU-001 – Realizar <i>Stemming</i> em um texto</b>	
<b>Objetivo:</b>	Realizar processo de <i>Stemming</i> sobre um texto informado
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário (identificado como Ator) e Algoritmo de <i>Stemming</i> (identificado como sistema)
<b>Cenário Principal:</b>	
1. Ator informa sobre qual texto deseja realizar o processo de <i>Stemming</i>	
2. Ator informa sobre qual língua (ex. Portuguesa) deseja que o sistema faça o <i>Stemming</i>	



3. Sistema processa texto utilizando dicionário da língua selecionada
4. Sistema analisa cada palavra contida no texto informado pelo Ator
5. Sistema aplica algoritmo de <i>Stemming</i> sobre cada palavra analisada
6. Sistema retorna texto modificado para o Ator, contendo somente os <i>stems</i> (radicais) das palavras contidas no texto originalmente informado
7. Caso de uso se encerra
<b>Cenários Alternativos:</b>
Sistema não consegue realizar o <i>Stemming</i> de uma palavra:
5.a. Sistema não consegue realizar <i>Stemming</i> sobre determinada palavra do texto informado, desta maneira ele retorna a palavra na sua forma original

### 8.5.1.2 Diagrama de Classes do Processo de *Stemming*



**Figura 8.1 – Diagrama de Classes do Processo de *Stemming***

A figura 8.1 mostra as principais classes envolvidas com o processo de *Stemming*. A classe “*PortuguesStemmer*” é a responsável por receber o texto informado pelo usuário. Nesta mesma classe é instanciado um novo objeto para a utilização da classe “*PortugueseAnalyzer*”.

A classe “*PortugueseAnalyzer*” por sua vez recebe os parâmetros repassados pela classe anterior e instancia uma nova classe, chamada de “*SnowballAnalyzer*”.

A nova classe instanciada, pertencente à *API Lucene*, é a responsável por chamar os métodos que atuarão sobre as informações passadas pelo usuário, retornando o texto somente com o resultado do *Stemming*.

A classe “*Analyzer*” é a classe pai das classes “*PortugueseAnalyzer*” e “*SnowballAnalyzer*” provendo alguns métodos para as mesmas.

### 8.5.1.3 Diagrama de Seqüência do processo de *Stemming*

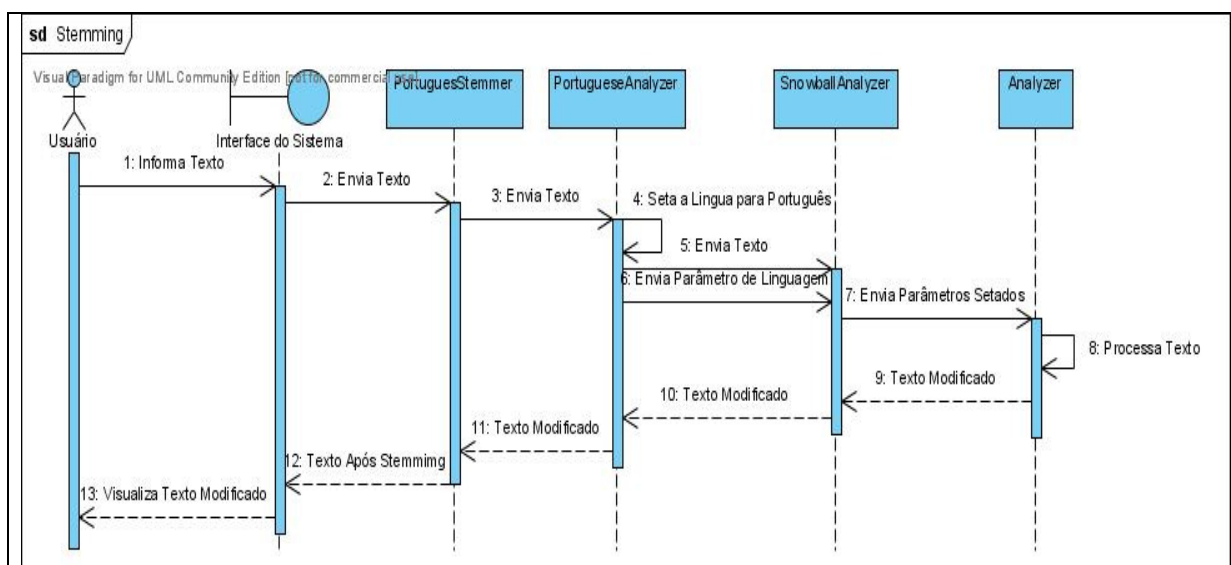


Figura 8.2 – Diagrama de Seqüência do processo de *Stemming*

### 8.5.2 Algoritmo de Remoção de *Stopwords*

Nesta seção é apresentada a modelagem realizada sobre o algoritmo de remoção de *Stopwords*. Este algoritmo que pertence a *API Lucene* tem a função de realizar a remoção de *Stopwords* de um dado texto.

#### 8.5.2.1 Realizar processo de Remoção de *Stopwords* de um texto

Nesta subseção é mostrado o principal caso de uso extraído através do desenvolvimento de um exemplo utilizando o algoritmo de remoção de *stopwords*.

<b>Objetivo:</b>	Realizar processo de remoção de <i>Stopwords</i> sobre um texto informado
<b>Pré-Condição:</b>	N/A
<b>Atores:</b>	Usuário (identificado como Ator) e Algoritmo de Remoção de <i>Stopwords</i> (identificado como sistema)
<b>Cenário Principal:</b>	
1. Ator informa sobre qual texto deseja realizar o processo de Remoção de <i>Stopwords</i>	
2. Sistema recebe as informações passadas pelo Ator para processamento	
3. Sistema processa texto utilizando a lista de <i>stopwords</i> disponíveis para execução	
4. Sistema analisa cada palavra contida no texto informado pelo Ator	
5. Se a palavra analisada estiver contida na lista de <i>stopwords</i> disponíveis, o sistema a remove. Caso contrário continua o processo	
6. Sistema retorna o texto modificado (sem as <i>stopwords</i> ) para o Ator	
7. Caso de uso se encerra	
<b>Cenários Alternativos:</b>	
Ator informa sua própria lista de <i>stopwords</i> :	
1.a. Ator informa juntamente com o texto uma lista própria de <i>stopwords</i> , o sistema passa a utilizar a mesma no restante do processamento	

### 8.5.2.2 Diagrama de Classes do Processo de Remoção de *Stopwords*

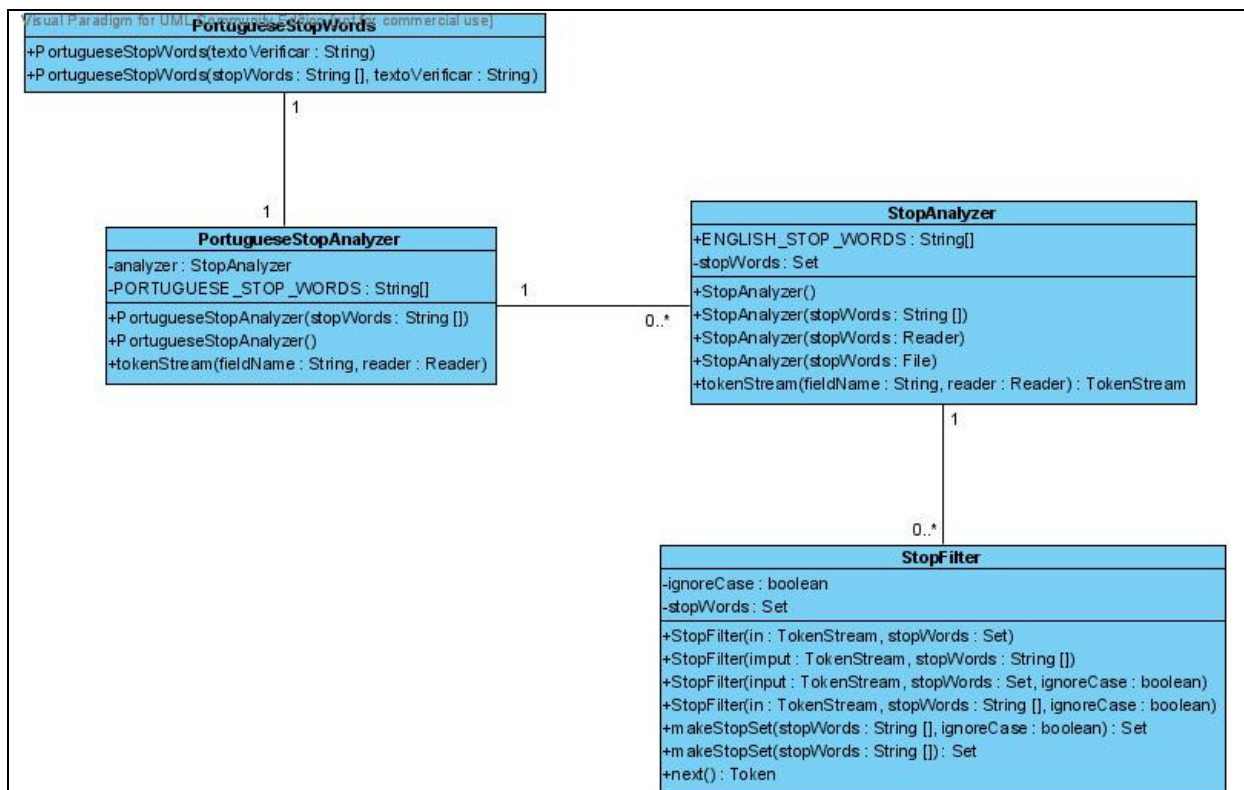


Figura 8.3 – Diagrama de Classes do Processo de Remoção de *Stopwords*

O diagrama 8.3 demonstra as principais classes envolvidas no processo de remoção de *Stopwords*.

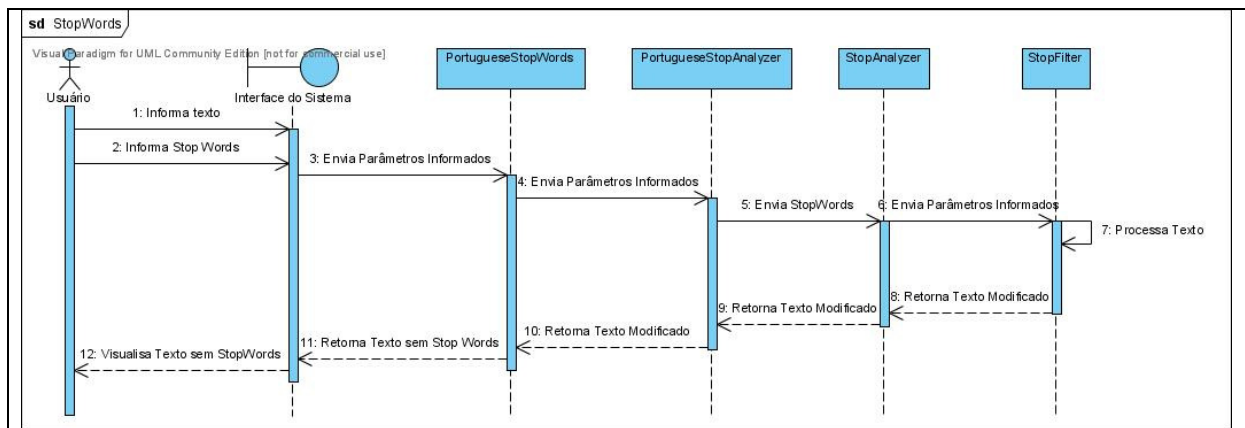
A classe “*PortugueseStopWords*” é a classe responsável por receber o texto ou as *stopwords* e o texto informados pelo usuário do sistema. Nesta classe é instanciado um novo objeto para a utilização da classe “*PortugueseStopAnalyzer*” passando exatamente os parâmetros informados pelo usuário.

Na classe “*PortugueseStopAnalyzer*”, é instanciado um novo objeto para a classe “*StopAnalyzer*”, classe pertencente a *API Lucene*, repassando o texto e as *stopwords*.

A classe “*StopAnalyzer*” instancia um novo objeto para outra classe da *API Lucene*, a classe “*StopFilter*”.

A classe “*StopFilter*” é a responsável por chamar os métodos que irão realizar o processo de Remoção das *Stopwords* do texto.

### 8.5.2.3 Diagrama de Seqüência do processo de Remoção de *Stopwords*



**Figura 8.4 – Diagrama de Seqüência do Processo de Remoção de *Stopwords***

## 9 MODELAGEM DO *FRAMEWORK* – DIAGRAMAS DE CLASSES

Neste capítulo é apresentada a modelagem do *framework* proposto para o corrente trabalho. Nesta modelagem são considerados apenas os diagramas de classes, pois as principais funcionalidades que estão aqui representadas se referem às ferramentas de PLN modeladas no capítulo anterior.

O papel de cada classe nos diagramas seguintes será explicado no próximo capítulo, intitulado Prototipação.

### 9.1 Pacote *search*

Nesta seção é apresentado o diagrama de classes do pacote *search*. Este pacote é responsável por implementar as funcionalidades disponíveis nas *API's* dos mecanismos de busca. A figura 9.1 mostra este diagrama.

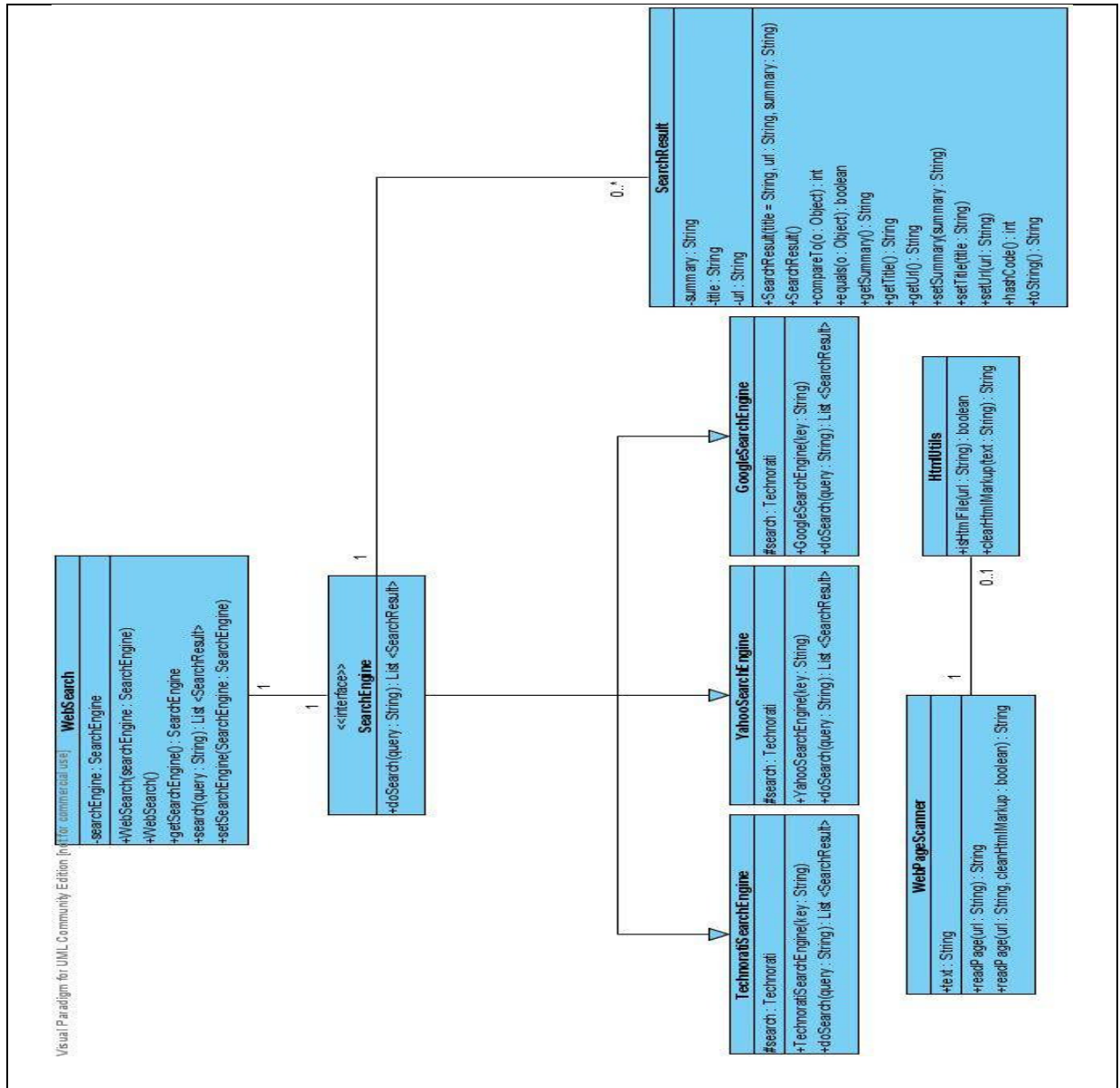
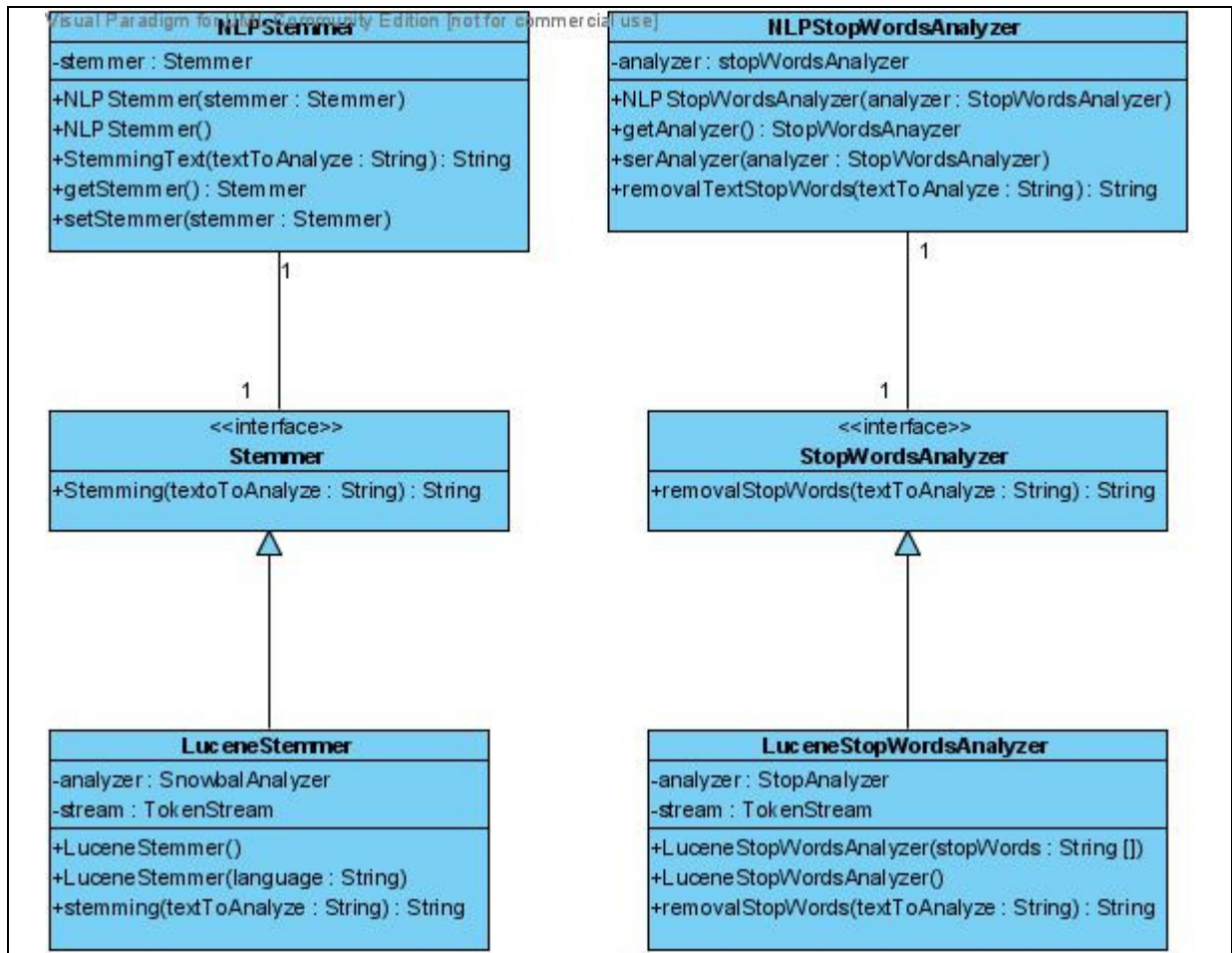


Figura 9.1 – Diagrama de Classes do Pacote *search*

## 9.2 Pacote *nlp*

Nesta seção é apresentado o diagrama de classes do pacote *nlp*, que implementa as ferramentas de PLN. A figura 9.2 mostra este diagrama.



**Figura 9.2 – Diagrama de Classes do Pacote *nlp***

### 9.3 Diagrama de Classes do Pacote *nlpSearch*

Nesta seção é apresentado o diagrama de classes do pacote *nlpSearch*. Neste pacote encontram-se as classes que integram os conceitos referentes a utilização das *API's* de mecanismos de buscas e as ferramentas de PLN.

O primeiro diagrama, figura 9.3, representa o conjunto de classes que foram utilizados para construir a classe *CorpusSearch*. Nesta classe são utilizadas funcionalidades das *API's* de mecanismos de busca para procurar páginas que contenham um termo informado, as funcionalidades da classe *WebPageScanner* para baixar o conteúdo das páginas encontradas e funcionalidades da classe *HtmlUtils* para verificar se cada página encontrada se refere a uma página no formato *HTML*.



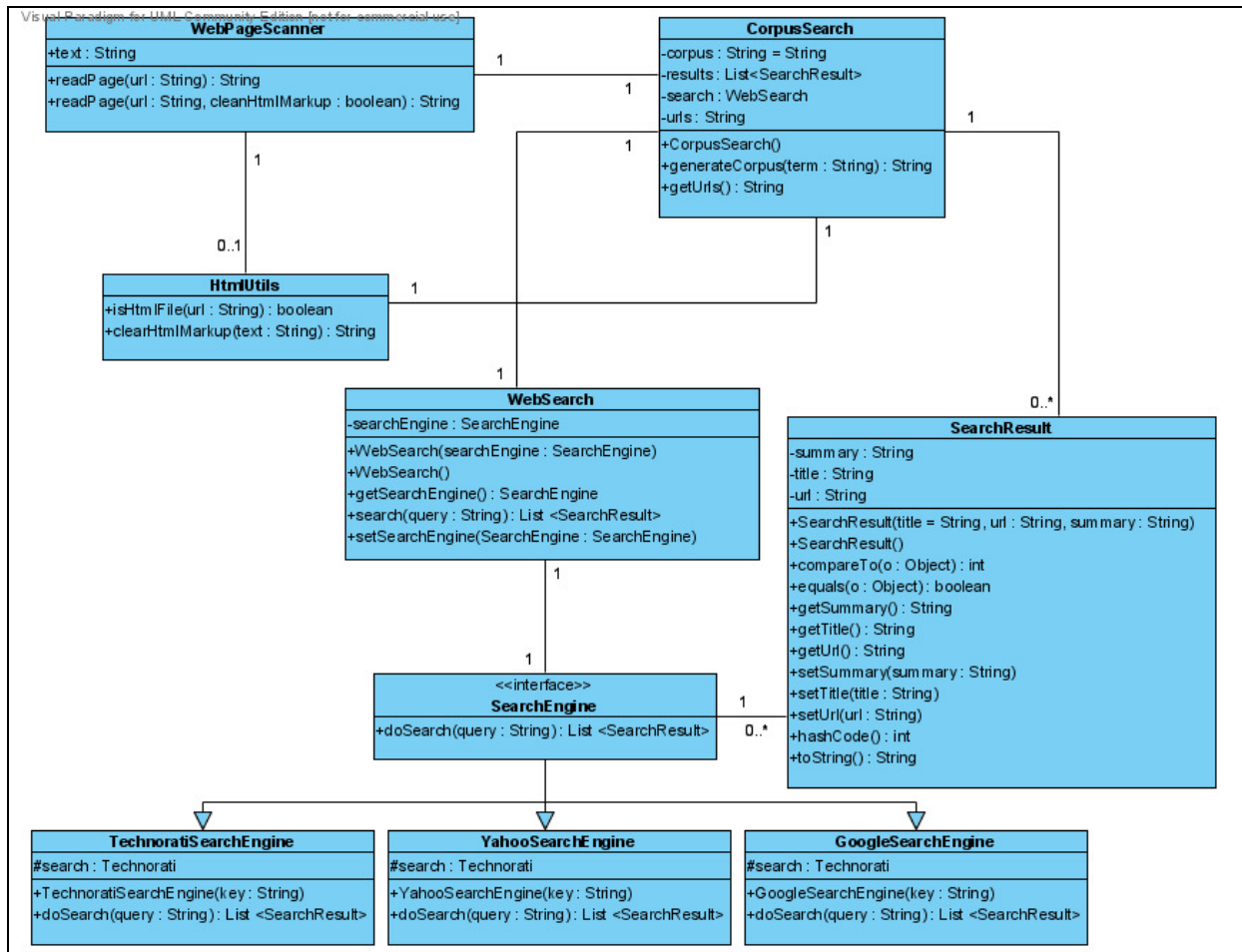


Figura 9.3 – Diagrama de Classes *CorpusSearch*

O diagrama seguinte, figura 9.4, representa o conjunto de classes utilizados para construir a classe *StemmingSearch*. O primeiro conjunto se refere às *API's*, enquanto o segundo se refere ao conjunto de funcionalidades referentes ao PLN. Assim como na classe anterior os mecanismos de busca servem para realizar as pesquisas na *web*, porém essas buscas serão feitas sobre termos que foram anteriormente processados, ou seja, esses termos sofreram alterações após ser realizado um processo de *stemming*.



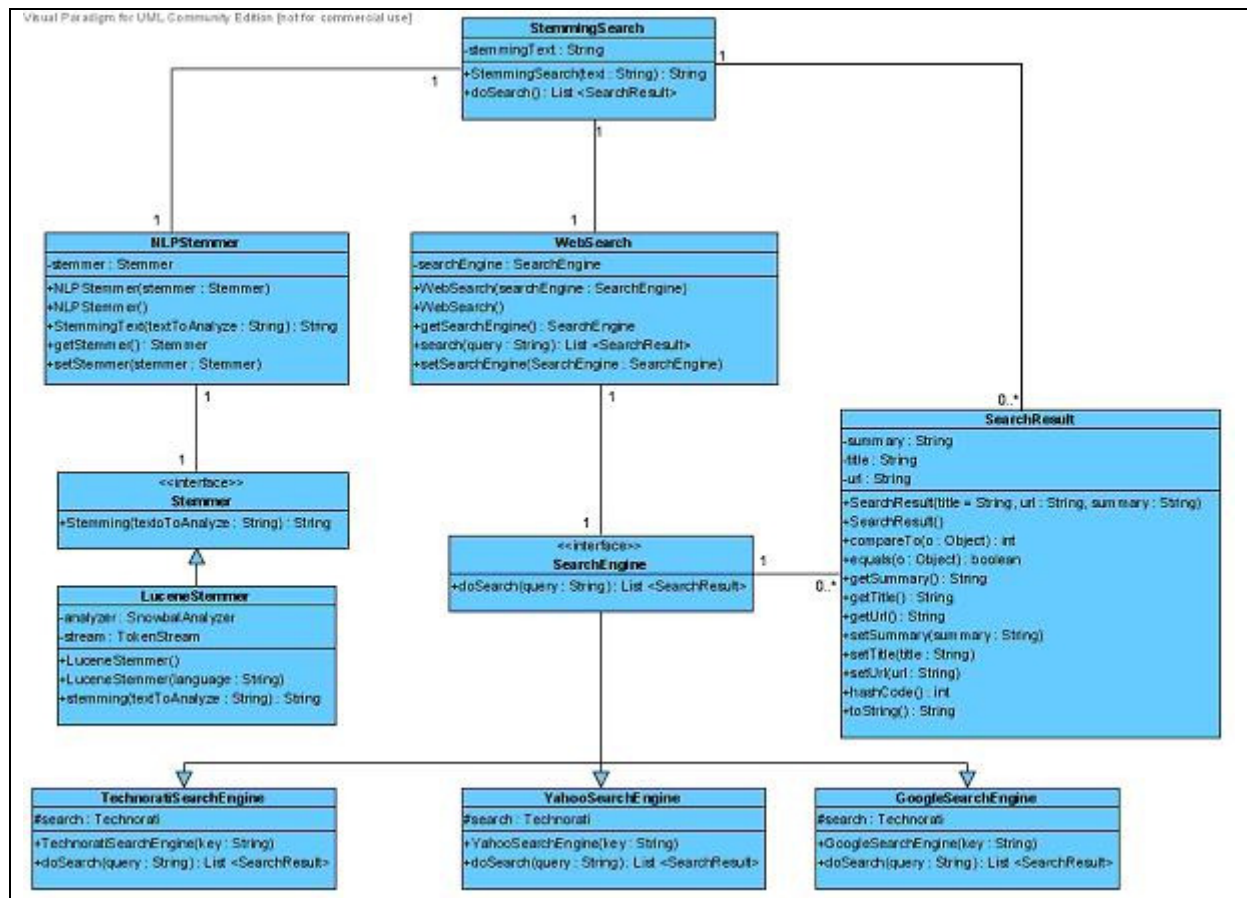


Figura 9.4 – Diagrama de Classes *StemmingSearch*

#### 9.4 Diagrama de Classes da Visão Geral do *Framework NLPSearch*

A figura 9.5 representa a visão geral do *Framework NLPSearch*, onde pode-se visualizar todas as classes que fazem parte dos pacotes que compõe o *framework* desenvolvido para o presente trabalho.

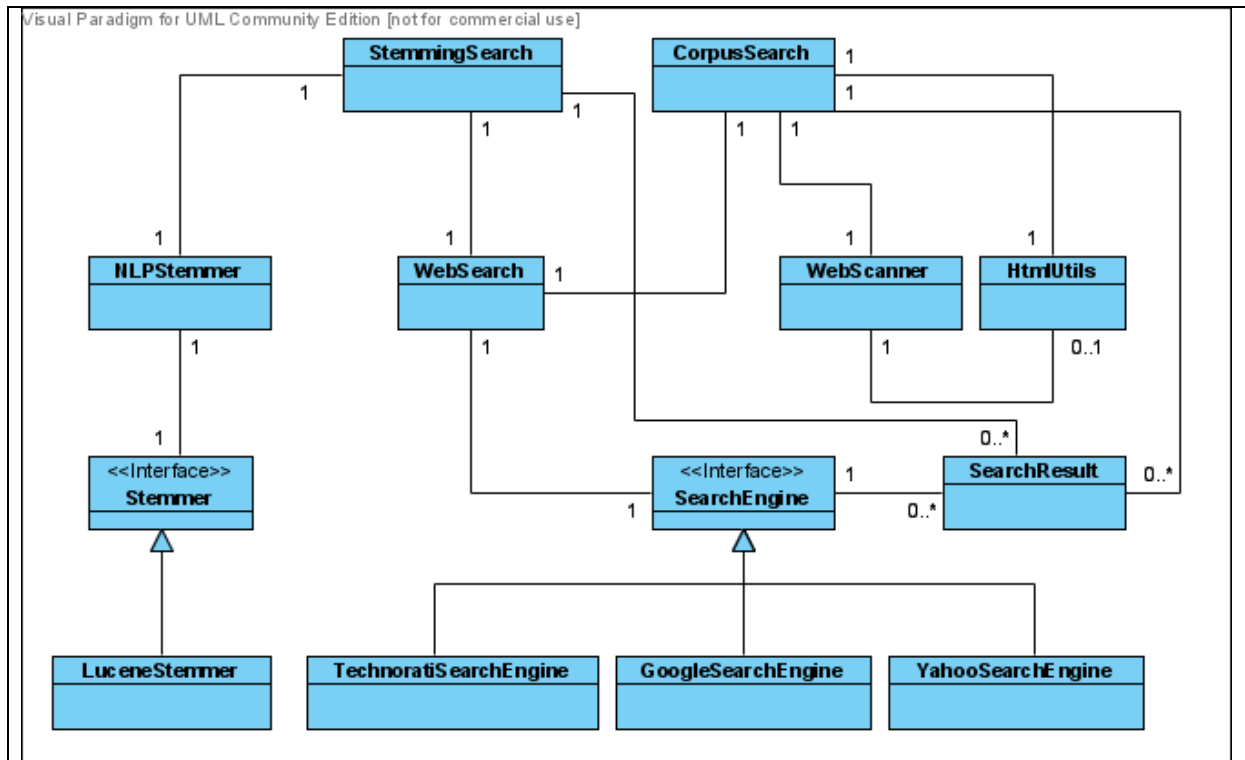


Figura 9.5 – Diagrama de Classes da Visão Geral do *Framework NLPSearch*

## 10 PROTOTIPAÇÃO

Neste capítulo são apresentados os protótipos desenvolvidos através da utilização da estrutura do *framework* discutido durante este trabalho. Os protótipos representam a utilização de algumas funcionalidades implementadas no *framework*.

O capítulo está dividido em duas principais seções, a primeira descreve os protótipos desenvolvidos, enquanto a segunda apresenta como o *framework* deve ser utilizado para gerar novas aplicações.

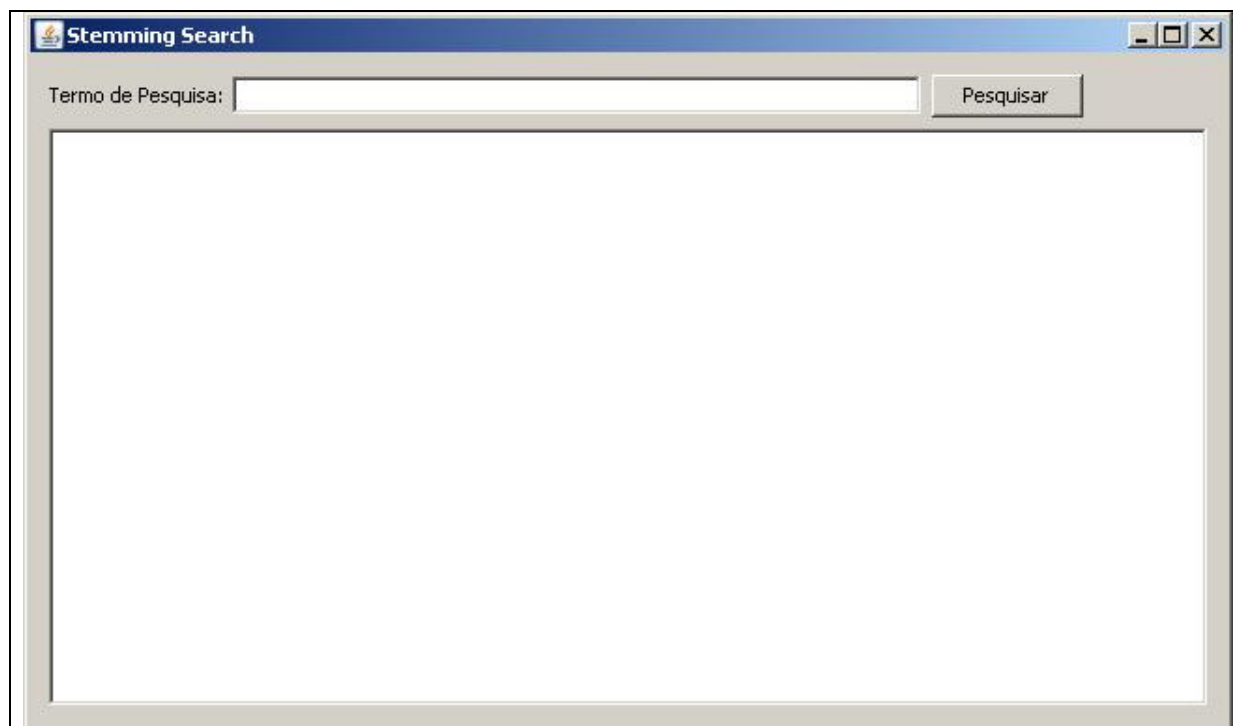
### 10.1 Funcionalidades Implementadas

Os protótipos que são apresentados neste capítulo implementam duas funcionalidades que integram as *API's* de mecanismos de busca e conceitos de PLN.

Os protótipos aqui apresentados estarão disponíveis na documentação eletrônica que será disponibilizada junto com o presente trabalho.

#### 10.1.1 Protótipo *Stemming Search*

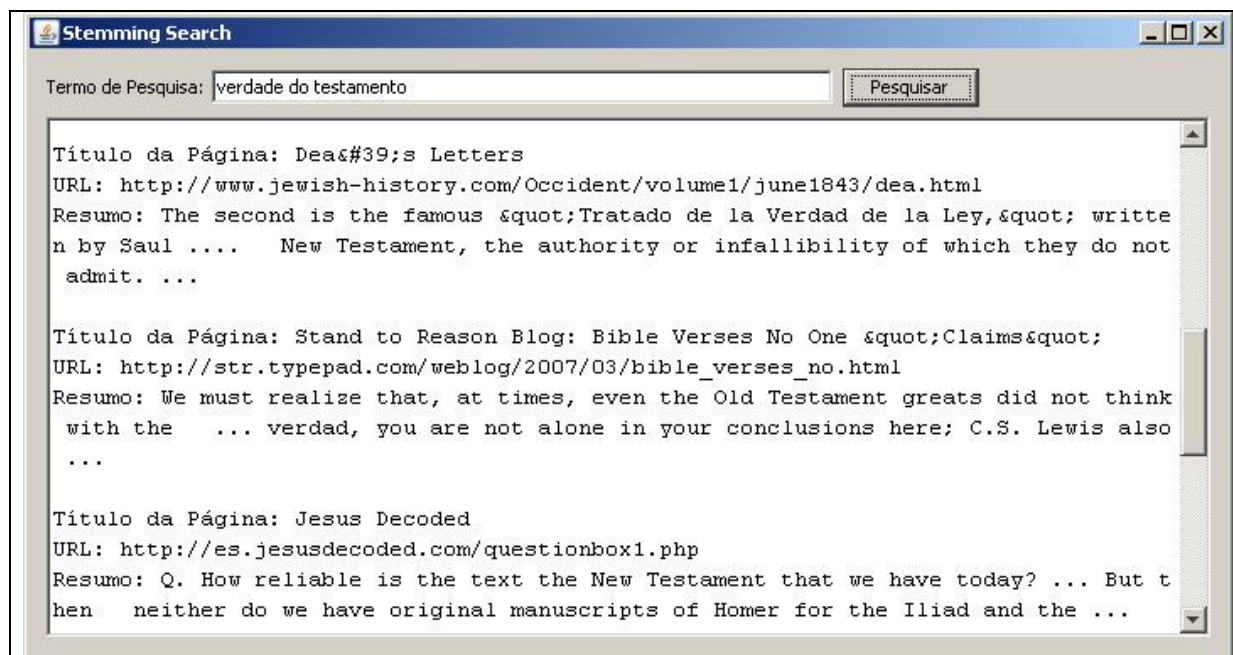
A figura 10.1 apresenta a tela da aplicação desenvolvida para integrar conceitos de *API's* e o processo de *Stemming*, conceito de PLN. Ela é composta basicamente por duas principais áreas, a área onde são informados os termos que serão pesquisados e a área onde serão apresentados os resultados. O botão “Pesquisar” é o responsável por realizar a pesquisa através da *API* do *Google*, implementada para este exemplo de aplicação.



**Figura 10.1 – Aplicação *Stemming Search***

A figura 10.2 apresenta o teste realizado com a aplicação desenvolvida, onde foi informado um termo de pesquisa e solicitado ao sistema que ele exiba páginas onde a palavra está presente. Neste caso é importante lembrar que as palavras digitadas no campo de “Termo de Pesquisa” irão passar por um processo de *stemming*, e após este processo serão enviadas para o mecanismo de buscas, onde será realizada a pesquisa. Na área inferior do protótipo pode – se observar os resultados da pesquisa, para cada página retornada são exibidas três informações: Título da Página, URL e Resumo.

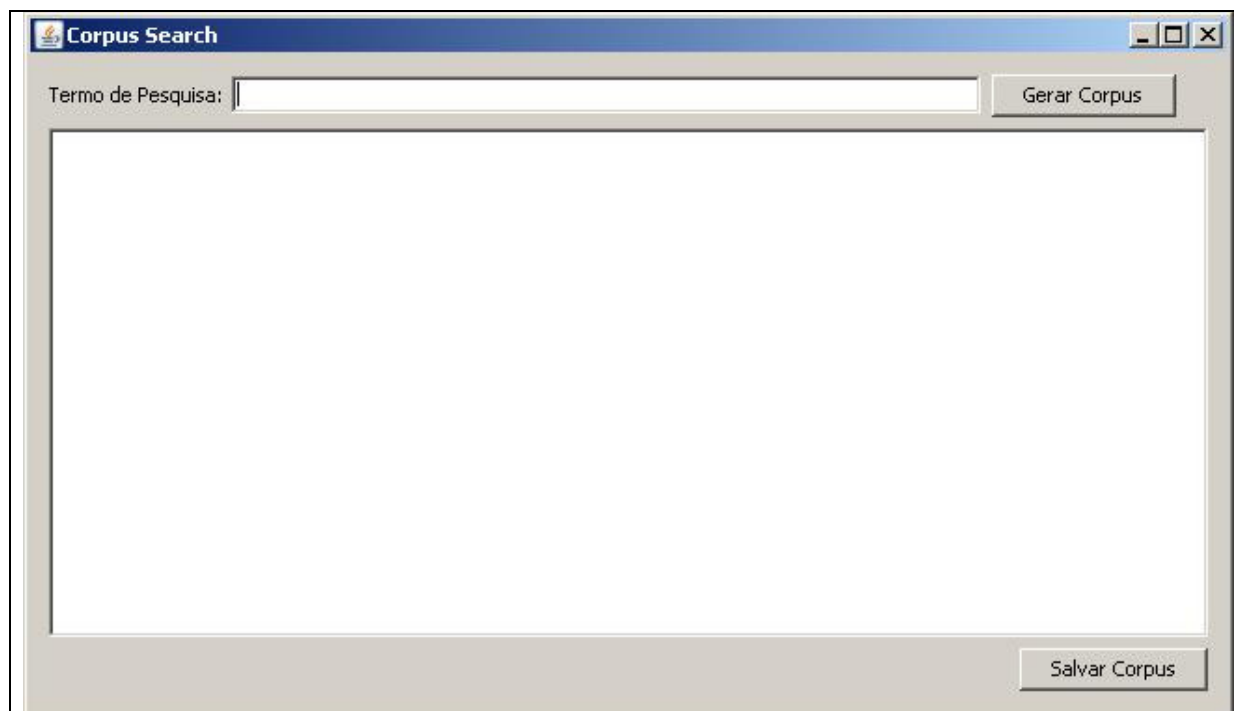
Pode-se notar ao observar estas palavras (palavras pesquisadas) que o processo de *stemming* removeu suas terminações, ou seja, a palavra “testamento” foi pesquisada como “testament” enquanto a palavra “verdade” foi pesquisada como “verdad”.



**Figura 10.2 – Resultado da Aplicação *Stemming Search***

### **10.1.2 Protótipo *Corpus Search***

A figura 10.3 apresenta a tela da aplicação *CorpusSearch*, desenvolvida para mostrar uma funcionalidade disponível no *framework*, ou seja, a construção de um *corpus* através da união de textos coletados em páginas da *web*. A aplicação é composta por duas áreas principais, a área de pesquisa de termos e a área onde será apresentado o *corpus* gerado. O botão “Gerar Corpus” é o responsável por realizar a pesquisa na *web* e montar os textos que irão compor o *corpus*, assim como na aplicação anterior o mecanismo de busca escolhido foi o do *Google*. O botão “Salvar Corpus” fornece a opção ao usuário de salvar o *corpus* em algum arquivo escolhido por ele.

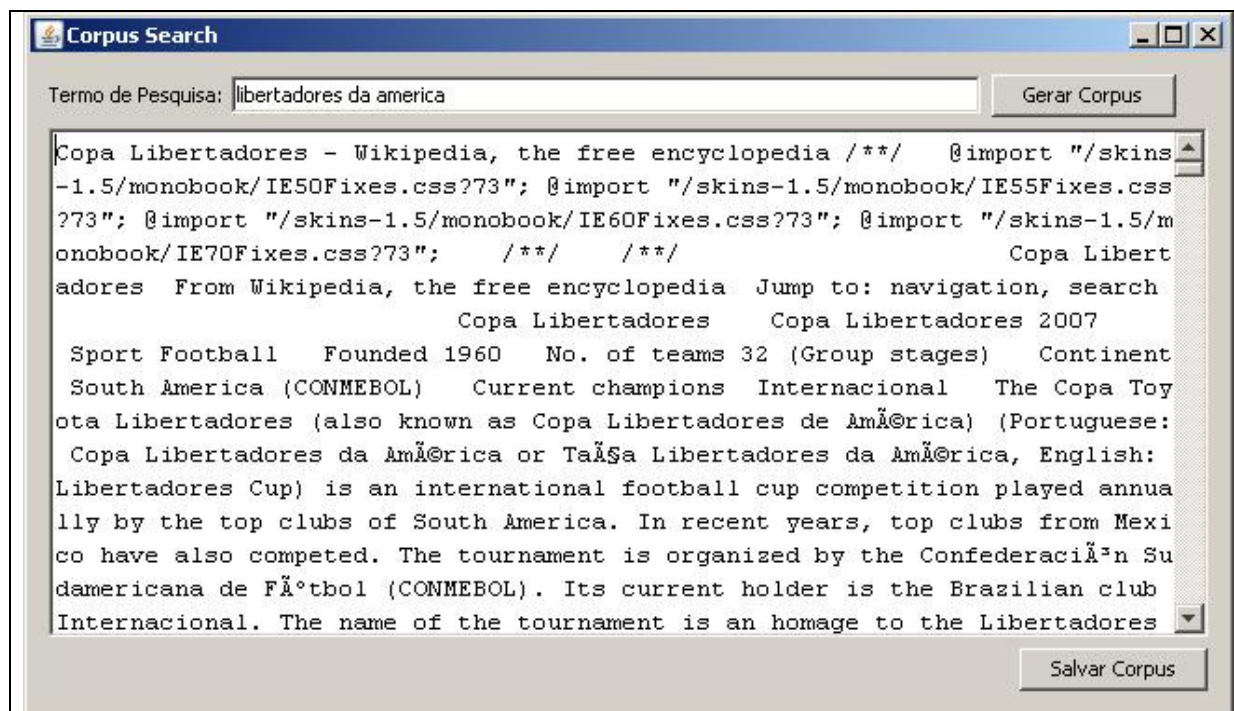


**Figura 10.3 – Aplicação *Corpus Search***

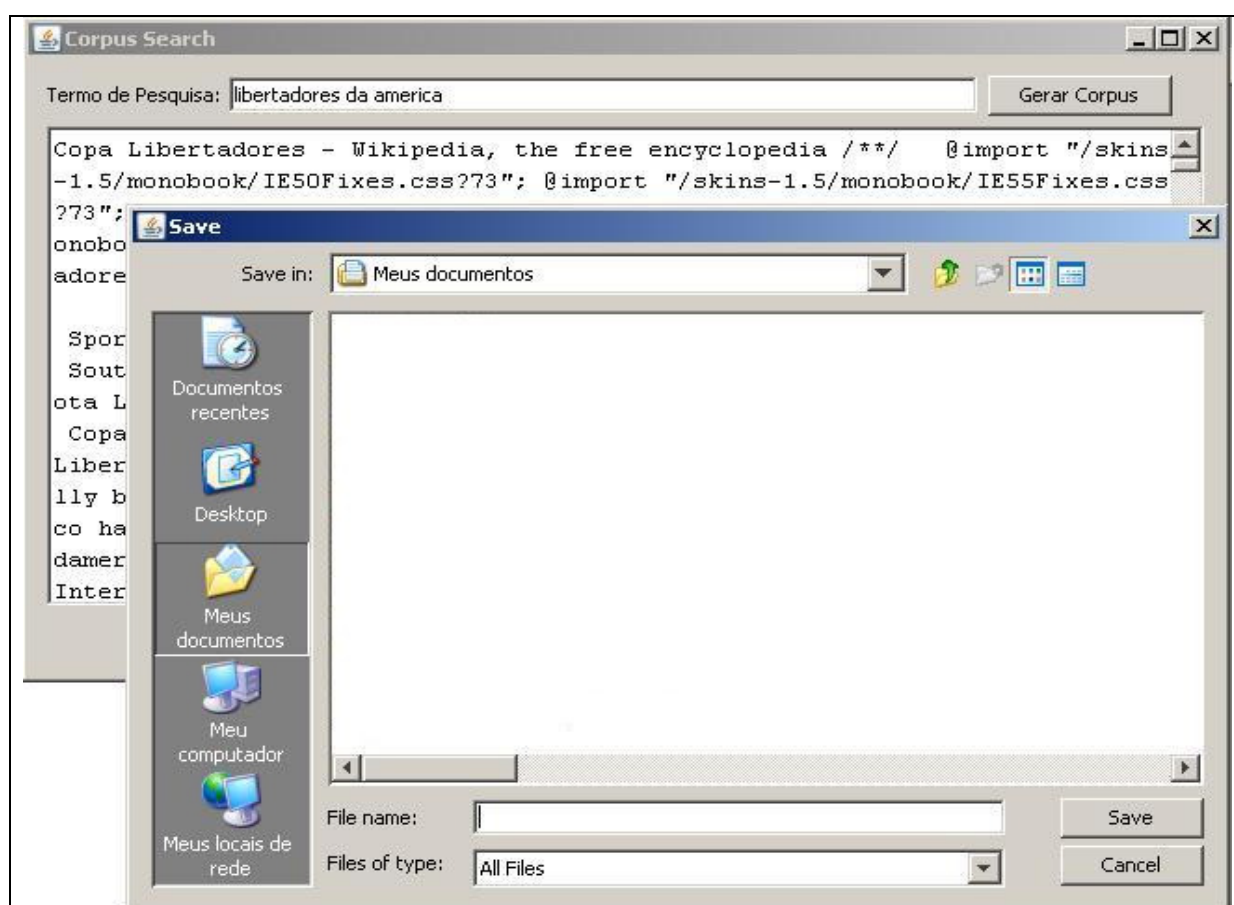
As próximas figuras apresentadas mostrarão os testes realizados com o protótipo de aplicação. A figura 10.4 apresenta o teste inicialmente realizado, onde foi informado um termo de pesquisa no campo “Termo de Pesquisa” e solicitado através do botão “Gerar Corpus” que o sistema buscasse na *web* por páginas que contenham o termo e gerasse o *corpus*. É importante ressaltar aqui que o processo para gerar o *corpus* passa por algumas etapas, entre elas a remoção de marcações de páginas *web*, como por exemplo, marcações que ficam entre “<?>”.

Na área inferior da tela é apresentado o *corpus* gerado, assim como o botão que permite ao usuário salvar em um arquivo o *corpus* gerado.

A figura 10.5 mostra o processo realizado para salvar o *corpus* em um arquivo, enquanto a figura 10.6 mostra o arquivo salvo.



**Figura 10.4 – Corpus gerado através da aplicação *Corpus Search***



**Figura 10.5 – Salvando o *Corpus* gerado**



```

corpus.txt - Bloco de notas
Arquivo Editar Formatar Exibir Ajuda
Copa Libertadores - wikipedia, the free encyclopedia /.../ @import "/skins-1.5/monobook/IE50Fixes.css?73"; @import
nacional. The name of the tournament is an homage to the Libertadores (Spanish word for Liberators), the main leade
a, Brazil, Chile, Colombia, Ecuador, Paraguay, Peru and Uruguay) would play in the Libertadores Cup, but in the 197
m. The top two teams from each group are then drawn in the second round, which consists of a two-legged knockout ti
se extra time or away goals to decide a tie that was level on aggregate. From 1960 to 1987, two-legged ties were de
in the second round, and would also have lost the final on away goals, making Boca Juniors the Champion of that trc
nners seven times including four in a row from 1972 to 1975. The competition all-time goalscorer is still striker A
er the years the competition has kept alive a healthy sport rivalry between the competing countries, especially bet
1964 Details Independiente Argentina 0 - 0 1 - 0 Aggregate 1 - 0 Nacional Uruguay Santos FC Brazil Colo-Colo CF
il Barcelona Sporting Ecuador 1972 Details Independiente Argentina 0 - 0 2 - 1 Aggregate 2 - 1 Universitario F
nte Argentina Guarani Brazil Year Final Semifinalists Winner Score Runner-up 1980 Details Nacional Uruguay
de Cali Colombia Cobreloa Chile River Plate Argentina 1988 Details Nacional Uruguay 0 - 1 3 - 0 Aggregate 3 -
gregate 1 - 1 (pen) 5 - 3 São Paulo FC Brazil AtlÁtico Junior Colombia Olimpia Paraguay 1995 Details GrÃmic
(pen)3 - 1 Cruz Azul Mexico Palmeiras Brazil Rosario Central Argentina 2002 Details Olimpia Paraguay 0 - 1 2
) 3 times Estudiantes L.P. (1968, 1969, and 1970) Nacional (1971, 1980, and 1988) Olimpia (1979, 1990, and 200
pscorers See Copa Libertadores de AmÃrica Topscorers [edit] External links Copa Toyota Libertadores at CONMEBOL.
#160;#160;Asia: AFC Ãe" Champions League #160;#160;#160;#160;#160;Africa: CAF Ãe" Champions League #160;#160;#160;
ion
Permanent linkCite this article In other languages @SU,'0'±0
(i=0;!x&&i
Ideo | Clubes | Fale conosco globoEsporte.com Futebol 2007 Libertadores _uacct = "UA-
"/ @import "/Portal/cda/estilo_css_cda/0,,3259,00.css"; /* globoesportev2_topo_Libertadores */ @import "/Portal/cda
ment.write(''); } } //-- OAS_PrintGLB(); OAS_AD('top'); OAS_PrintGLB(); OAS_AD('Frame1');
n false; } } } function menuToogle(obj) { obj = obj.parentNode; if(obj.className.indexOf('nive1A')-1) { if(obj.c
ndÃrio Ãltimas NotÃcias Surfe Capa de Surfe Blog do Teco Galerias de Fotos Rico Surfe Supersurf
ricano) Capa de NFL A Liga Blog da NFL Entenda o Jogo Galerias de Fotos Times Tabela de Jogos Ãlti
favorito Boca e ficam perto da decisÃo
A nova versÃo da febre da internet AtÃ o papa jÃ; jogou! Clique
o 4 PG J V E D GP GC SG 01 VÃlez Sarsfield 11 6 3 2 1 6 3 3 02 Nacional 10 6 3 1 2 9 5 4 03 Internac
Ver o regulamento &raquo; Ver a Tabela de Jogos &raquo; Fase Anterior OITAVAS PrÃxima Fase
val de BritoAssunÃÃo Ver o regulamento &raquo; Ver a Tabela de Jogos &raquo; Fase Anterior
Ver o regulamento &raquo; Ver a Tabela de Jogos &raquo; Fase Anterior FINAL Data
VÃ-deos Ã»
Libertadores 2007 Ãltimos Jogos Quinta, 31/05/2007 CÃcuta Deportivo 3 x 1 Boca Juniors C
3"; @import "/skins-1.5/monobook/IE60Fixes.css?73"; @import "/skins-1.5/monobook/IE70Fixes.css?73";
s 5 ver tambÃom // [editar] ClassificaÃÃo A forma de classificaÃÃo para a competiÃÃo Ã geralmente bas
undo colocado, o 2Ã melhor joga contra o penÃltimo dos segundos colocados, e assim por diante. A disputa acontece
rneio Apertura: 1 #160;MÃxico: #160;Campeonato Mexicano: 2 #160;Torneio Apertura e Torneio ClausuraÃ: 1 &

```

Figura 10.6 – *Corpus* salvo

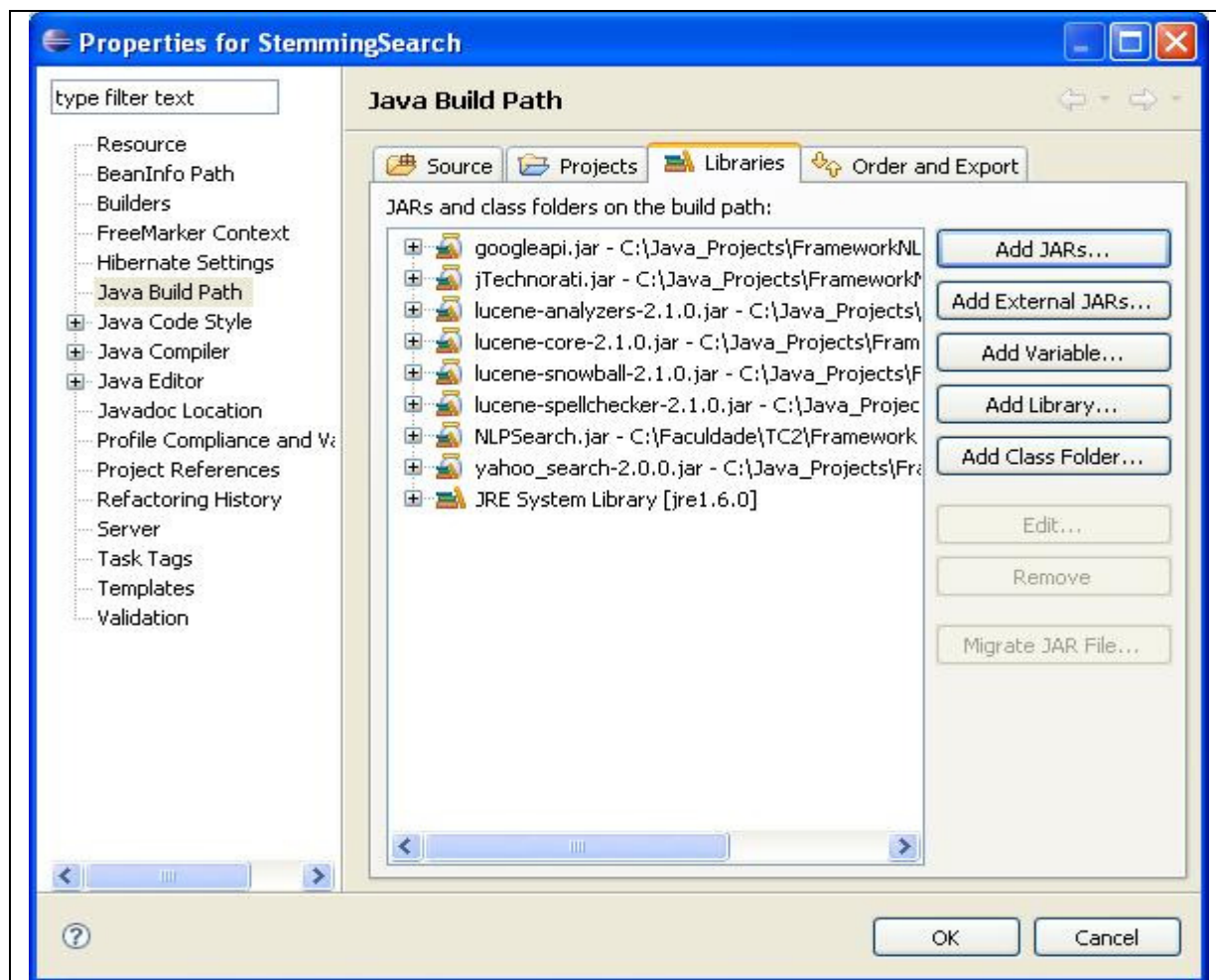
## 10.2 Como utilizar o *framework* para desenvolver aplicações

Nesta seção serão apresentados os principais requisitos para utilizar a estrutura do *Framework NLPSearch* para a construção de aplicações.

### 10.2.1 Estrutura necessária

Para utilizar as funcionalidades implementadas no *Framework* é necessário que o usuário atenda a seguinte estrutura de bibliotecas, apresentada na figura 10.7, na sua aplicação. E também atenda a algumas premissas básicas para realizar pesquisas na *web*.





**Figura 10.7 – Estrutura de diretórios**

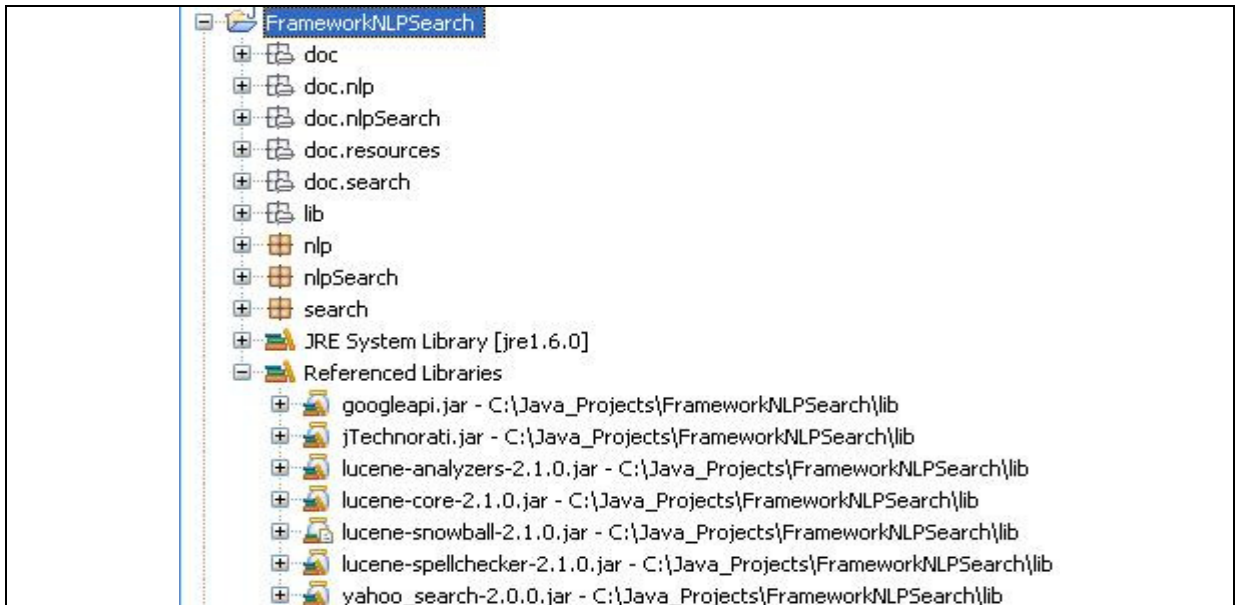
Como pode ser observado na figura anterior, o *Framework NLPSearch* também está inserido como sendo uma biblioteca da aplicação, através desta biblioteca, chamada *NLPSearch* o usuário poderá utilizar suas principais funcionalidades.

Para que o usuário do *framework* consiga construir e testar suas aplicações é necessária uma conexão com a internet, pois as aplicações desenvolvidas visam integrar as *API's* de mecanismos de busca (utilizam a internet para buscar informações) com ferramentas de PLN.

### 10.2.2 Estrutura de pacotes do *Framework NLPSearch*

A figura 10.8 mostra a composição do *framework*. Nos cinco primeiros pacotes estão armazenados os arquivos que contém a documentação do *framework*, documentação feita através da utilização do *javadoc*. No pacote “*nlp*” encontram-se as classes que implementam as funcionalidades de PLN. No pacote seguinte chamado de “*nlpSearch*” encontram-se as

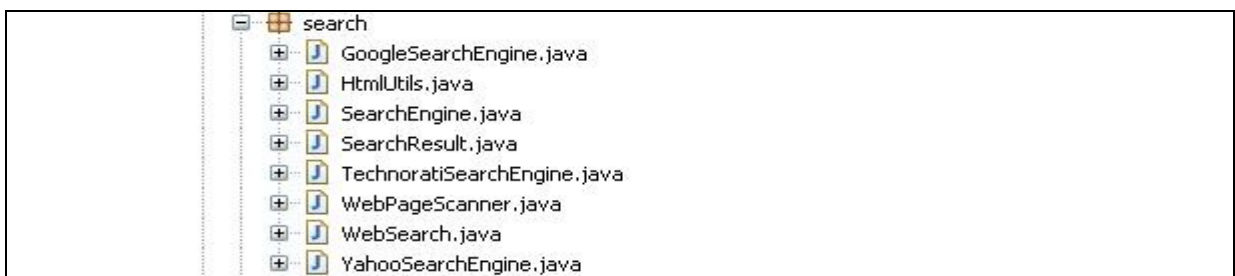
classes que integram as *API's* de mecanismos de busca e as ferramentas de PLN. No pacote “*search*” encontram-se as classes relacionadas com os mecanismos de busca e pesquisas na *web*.



**Figura 10.8 – Estrutura do *Framework NLPSearch***

### 10.2.3 Estrutura de Classes do *Framework NLPSearch*

Cada pacote apresentado anteriormente é composto por um conjunto de classes que contemplam algumas funcionalidades, estas funcionalidades serão apresentadas a seguir.



**Figura 10.9 – Estrutura de classes do pacote *search***

O pacote apresentado na figura 10.9 possui classes que realizam pesquisas na *web*, verificam tipos de arquivos, baixam páginas e guardam resultados de pesquisa. As classes “*GoogleSearchEngine*” (*hot-spot*), “*YahooSearchEngine*” (*hot-spot*), e “*TechnoratiSearchEngine*” (*hot-spot*), são as responsáveis por realizar buscas na *web*. A classe “*SearchEngine*” é a interface para utilização da classe anteriormente citada. A classe “*WebSearch*” é a responsável por instanciar as classes dos mecanismos de busca. A classe

“*HtmlUtils*” possui algumas funcionalidades que podem ser aplicadas sobre arquivos *HTML*. Uma delas verifica se a URL encontrada é realmente pertencente a um arquivo *HTML*, outra retira as marcações *HTML* de uma página, deixando somente o que não faz parte destas marcações. A classe “*WebPageScanner*” é responsável por baixar todo o conteúdo de uma página *web*.



**Figura 10.10 - Estrutura de classes do pacote *nlp***

O pacote apresentado na figura 10.10 apresenta as classes implementadas para realizarem funcionalidades relacionadas ao Processamento da Linguagem Natural. A classe “*LuceneStemmer*” (*hot-spot*) realiza o processo de *Stemming* sobre as palavras para ela repassadas. Este *Stemmer* pode realizar o processo em diversas línguas, dentre elas está o português. A classe “*Stemmer*” é utilizada como uma interface para a utilização da classe anteriormente citada. A classe “*NLPStemmer*” é a responsável por instanciar a classe “*LuceneStemmer*”.

A classe “*LuceneStopWordsAnalyzer*” (*hot-spot*) realiza o processo de remoção de *stop words* de textos. Estas *stop words* podem ser repassadas para a classe, ou caso contrário será utilizada uma lista de *stop words* na língua inglesa. A classe “*StopWordsAnalyzer*” é utilizada como uma interface para a utilização da classe anteriormente citada. A classe “*NLPStopWordsAnalyzer*” é a responsável por instanciar a classe “*LuceneStopWordsAnalyzer*”.



**Figura 10.11 - Estrutura de classes do pacote *nlpSearch***

O pacote apresentado na figura 10.11 contém as classes que integram as *API's* de mecanismos de busca com as ferramentas de PLN. A classe “*CorpusSearch*” foi utilizada para implementar o protótipo *Corpus Search*, enquanto a classe “*StemmingSearch*” foi utilizada para desenvolver o protótipo *Stemming Search*.

### 10.2.4 Documentação do *Framework NLPSearch*

A documentação sobre a utilização das classes que compõem o *framework* está disponível no formato *javadoc*. A figura 10.12 apresenta esta documentação.

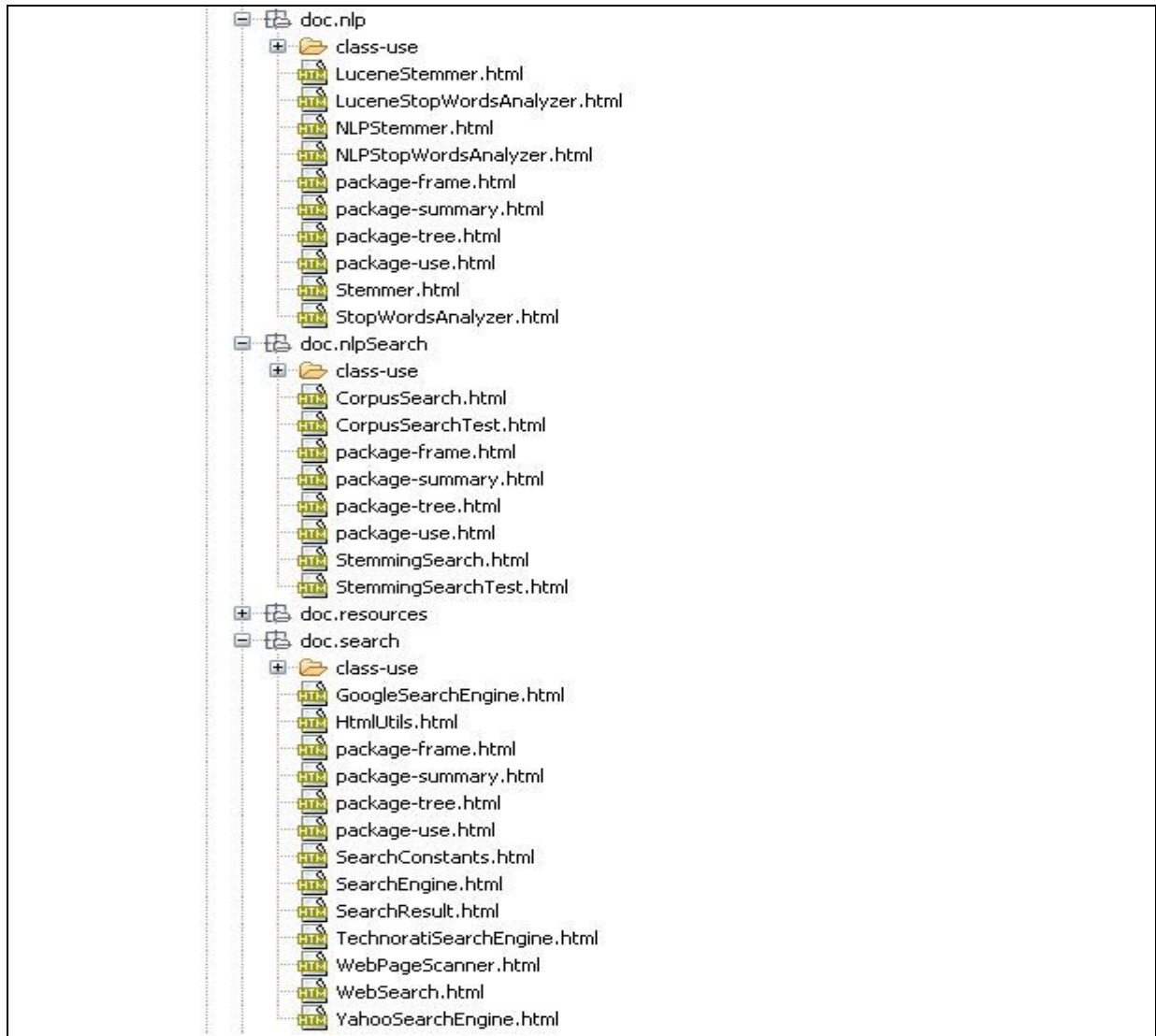


Figura 10.12 – Documentação do *Framework NLPSearch*

## CONSIDERAÇÕES FINAIS

O objetivo deste trabalho é a construção de um *framework* para desenvolvimento de aplicações para PLN que utilizem a *web* como *corpus*. Dessa forma foram apresentados diversos conceitos sobre orientação a objetos, *UML* principalmente para dar uma base concreta para o desenvolvimento do *framework* proposto.

O estudo das *API's* dos principais mecanismos de busca na *web*, como *Google*, *Yahoo* e *Technorati* viabilizou o entendimento de como as *API's* disponibilizadas gratuitamente por estes mecanismos podem auxiliar no decorrer do desenvolvimento do trabalho proposto. As três *API's* propostas para estudo geram formulários em *XML* que são disponibilizados como resultado para as pesquisas realizadas, facilitando dessa forma sua utilização, pois a linguagem *XML* está sendo utilizada largamente hoje em dia pela comunidade de *software*, tornando a compreensão dos resultados uma tarefa mais simples.

O estudo sobre o Processamento da Linguagem Natural focou-se mais nas categorias das ferramentas que utilizam seus recursos. Procurou-se mostrar quais suas principais funcionalidades, recursos e problemas encontrados. Foi extremamente importante o estudo destas categorias de ferramentas, pois o *framework* proposto tem a função de criar aplicações que utilizem tais recursos.

Quando se deseja construir um *framework* deve-se ter em mente qual será seu propósito, dessa forma, é importante conhecer exemplos dos *softwares* que serão desenvolvidos a partir das estruturas desse *framework*. Sendo assim, foram estudadas quatro ferramentas que utilizam recursos de PLN, e estão enquadradas nas categorias apresentadas.

O *framework* desenvolvido no presente trabalho teve a intenção de unir os conceitos referentes às *API's* de mecanismos de busca (*Google*, *Yahoo* e *Technorati*) e os conceitos referentes às ferramentas de PLN. Para tanto cada ferramenta foi modelada com a intenção de extrair suas principais funcionalidades, para que estas fossem implementadas no trabalho

proposto. Para cada *API* foram desenvolvidas funcionalidades para buscar as informações por ela disponibilizadas na *web*, como por exemplo, o “Título de Página”, “URL” e “Resumos”. A partir destas informações tornou-se possível criar mecanismos que realizassem processos de PLN e que utilizassem a *web* para realizar pesquisas, e construir conjuntos de textos.

Com base no *Framework NLPSearch*, produto deste trabalho, foram desenvolvidas pequenas aplicações (protótipos) que visavam demonstrar a utilização das funcionalidades implementadas. Neste mesmo contexto foram apresentados os principais requisitos e formas para se utilizar da arquitetura do *framework* para a criação de novas aplicações.

Todo o processo de desenvolvimento do *framework* foi documentado em formato de *javadoc*, e estará disponível no material eletrônico que será entregue juntamente com este trabalho. As aplicações desenvolvidas para testar a veracidade do *framework* também estarão disponíveis neste mesmo material.

Para trabalhos futuros ficam idéias referentes ao desenvolvimento de mais funcionalidades dentro da estrutura do *framework*, como por exemplo, funcionalidades para realizar análise sintática e morfológica de textos extraídos da *web*, contagem de palavras, aplicação de sinônimos, etc. Com base na estrutura desenvolvida podem ser geradas diversas aplicações que utilizem a arquitetura do *framework* como base.

## REFERÊNCIAS BIBLIOGRÁFICAS

- ABRAHÃO, Paulo Ricardo Carneiro. **Modelagem de Implementação de um Léxico Semântico para o Português**. Porto Alegre: 1997. 118p. Dissertação (Mestrado em Informática) – Instituto de Informática, PUCRS, 1997. Disponível em: <<http://www.inf.pucrs.br/~linatural/publicacoes/PRicardoD1.zip>>. Acesso em: 20 ago. 2006.
- BAPTISTA, Jorge. **Estabelecimento e Formalização de Classes de Nomes Compostos**. Lisboa: 1994. 242p. Dissertação (Mestrado em Lingüística Portuguesa Descritiva) – Faculdade de Letras, Universidade de Lisboa, 1994. Disponível em: <<http://w3.ualg.pt/~jbaptis/download/JBaptista1994.htm>>. Acesso em: 17 out. 2006.
- BEZERRA, Eduardo. **Princípios de Análise e Projeto de Sistemas com UML**. Rio de Janeiro: Elsevier, 2002. 286 p.
- BIBER, Douglas; CONRAD, Susan; REPPEN, Randi. **Corpus Linguistics: Investigating Language Structure and Use** (Cambridge Approaches to Linguistics). Cambridge - USA: Cambridge University Press, 1998. 320p.
- BOOCH, Grady; RUMBAUGH, James; JACOBSON, Ivar. **UML, Guia do Usuário**. Rio de Janeiro: Campus, 2000. 472p.
- BRAGA, Rosana Teresinha Vaccare. **Um Processo para Construção e Instanciação de Frameworks Baseados em uma Linguagem de Padrões para um Domínio Específico**. São Paulo: 2003. 44p. Tese (Doutorado na Área de Ciências da Computação e Matemática Computacional) – Instituto de Ciências Matemáticas e de Computação, USP, 2003. Disponível em: <<http://www.teses.usp.br/download.php/teses/disponiveis/55/55134/tde-12052003-102000/publico/TeseFinalCorrigida.pdf>>. Acesso em: 17 out. 2006.
- BUSCHMAN, Frank et al. **Pattern-Oriented Software Architecture: A System of Patterns**. Chichester – ENG: John Wiley and Sons Ltd, 1996. 466p.
- CHAVES, Marcirio Silveira. **Um Estudo e Apreciação sobre Algoritmos de Stemming para a Língua Portuguesa**. In: Jornadas Iberoamericanas de Informática, 4., 2003a, Cartagena de Índias-Colômbia. **Anais...** Cartagena de Índias-Colômbia: 2003, s/p. Disponível em: <[http://xldb.di.fc.ul.pt/~mchaves/pg\\_portugues/public/stemming.pdf](http://xldb.di.fc.ul.pt/~mchaves/pg_portugues/public/stemming.pdf)>. Acesso em: 17 nov. 2006.
- CHAVES, Marcirio Silveira. **Uma Introdução à Lingüística Computacional**. In: Faculdade de Filologia da Universidade da Corunha – Espanha. 2003b, Corunha –ESP. 85p. Disponível em: <[http://xldb.di.fc.ul.pt/~mchaves/pg\\_portugues/master/Apresenta/Intro\\_LC\\_4pp.pdf](http://xldb.di.fc.ul.pt/~mchaves/pg_portugues/master/Apresenta/Intro_LC_4pp.pdf)>. Acesso em: 15 nov. 2006.
- COELHO, Jorge César B.; COLLOVINI, Sandra; VIEIRA, Renata. **Estudo de Corpus para Classificação de Expressões Anafóricas da Língua Portuguesa**. In: Workshop em Tecnologia da Informação e da Linguagem Humana (TIL 2005), 3., 2005, São Leopoldo. **Anais...** São Leopoldo: SBC, 2005. Disponível em:

<[http://www.inf.unisinos.br/~renata/laboratorio/publicacoes/til\\_scr\\_05.pdf](http://www.inf.unisinos.br/~renata/laboratorio/publicacoes/til_scr_05.pdf)>. Acesso em: 17 nov. 2006.

COLEMAN, Derek. **Desenvolvimento Orientado a Objetos: O Método Fusion**. Rio de Janeiro: Campus, 1996. 389p.

D'SOUZA, Desmond Francis; WILLS, Alan Cameron. **Objects, Components, and Frameworks with UML: The Catalysis Approach**. . New York - USA: Addison-Wesley Publishing Company, 1998. 785p.

DIAS, Maria Abadia L.; MALHEIROS, Marcelo de Gomensoro. **Automatic Extraction of Keywords for the Portuguese Language**. In: Workshop on Computational Processing of Written and Spoken Language - PROPOR'2006, 7., 2006, Itatiaia. **Anais...** Itatiaia: PROPOR'2006, 2006. p.204-207. Disponível em: <<http://ensino.univates.br/~mald/artigo-wsl.pdf>>. Acesso em: 17 nov. 2006.

DIAS, Maria Abadia L.; MALHEIROS, Marcelo de Gomensoro. **Estudo de Técnicas de Radicalização para a Língua Portuguesa**. In: Workshop de Engenharia e Tecnologia, 1., 2005, Lageado. **Anais...** Lageado: Univates Editora, 2005. Disponível em: <<http://scholar.google.com/url?sa=U&q=http://ensino.univates.br/~mald/artigo-wet.pdf>>. Acesso em: 17nov. 2006.

FAYAD, Mohamed E.; JOHNSON, Ralph E. **Domain-Specific Application Frameworks: Frameworks Experience by Industry**. Chichester – ENG: John Wiley and Sons Ltd, 1999. 682p.

FAYAD, Mohamed E.; Schimdt, Douglas C. **Object-Oriented Application Frameworks: Special Issue on Object-Oriented Application Frameworks**. Communications of the ACM, New York - USA, v.40, n.10, p.32-38, Oct. 1997. Disponível em: <<http://www.cs.wustl.edu/%7Eeschmidt/CACM-frameworks.html>>. Acesso em: 6 nov. 2006.

FLIP, **Flip On-line**. Disponível em: <<http://www.flip.pt/FLiPOnline/tabid/96/Default.aspx>>. Acesso em 15 nov. 2006.

FOWLER, Scott; SCOTT, Kendall. **UML Distilled: A Brief Guide to the Standard Object Modeling Language**. New York, USA: Addison-Wesley Publishing Company, 2000. 185p.

GAMMA, Erich et al. **Padrões de Projeto: Soluções Reutilizáveis de Software Orientado a Objetos**. Porto Alegre: Bookman, 2000. 364p.

GASPERIN, Caroline; GOULART, Rodrigo; VIEIRA, Renata. **Uma Ferramenta para Resolução Automática de Correferência**. In: IV Encontro Nacional de Inteligência Artificial – ENIA. XXII Congresso Nacional da Sociedade Brasileira de Computação. **Anais...** Campinas: SBC, 2003. p.163-172. Disponível em: <<http://www.inf.unisinos.br/~renata/laboratorio/publicacoes/art1.pdf>>. Acesso em: 14 nov. 2006.

GERBER, Luciano David. **Uma Linguagem de Padrões para Desenvolvimento de Sistemas de Apoio à Decisão Baseado em Frameworks**. Porto Alegre: 1999. 144p. Dissertação (Mestrado em Informática) – Faculdade de Informaica, PUCRS, 1999.

GONZALEZ, Marco Antonio I. **Termos e Relacionamentos em Evidência na Recuperação de Informação**. Porto Alegre: 2005. 182p. Tese (Doutorado em Ciência da Computação) – Programa de Pós-Graduação em Computação, UFRGS, 2005. Disponível em: <<http://www.inf.pucrs.br/~gonzalez/tr+/tesemarco.pdf>>. Acesso em: 17 nov. 2006.



- GOOGLE. **Google Soap Search API (beta)**. Disponível em: <[http://www.google.com/apis/reference.html#2\\_5](http://www.google.com/apis/reference.html#2_5)>. Acesso em: 29 ago. 2006.
- GUEDES, Gilleanes T. A. **UML 2: Guia de Consulta Rápida**. São Paulo: Novatec Editora Ltda, 2005. 109p.
- HARMON, Paul; WATSON, Mark. **Understanding UML: the developer's guide: with a Web-based application in Java**. San Francisco-USA: Morgan Kaufmann Publishers, Inc, 1998. 367p.
- JACKSON, Peter; MOULINIER, Isabelle. **Natural Language Processing for Online Applications: Text Retrieval, Extraction and Categorization**. Amsterdam/Philadelphia: John Benjamins Publishing Company, 2002. 226p.
- JOHNSON, Ralph E. **Documenting Frameworks Using Patterns**: Proceedings of the Conference on Object-Oriented Programming Systems Languages and Applications. SIGPLAN Notices, New York-USA, v.27, n.10, p.63-76, oct. 1992. Disponível em: <<http://citeseer.ist.psu.edu/johnson92documenting.html>>. Acesso em: 7 nov. 2006.
- JOHNSON, Ralph E. **Documenting Frameworks**. Framework Digest, New York – USA, v.1, n.13, oct. 1994. Disponível em: <<ftp://st.cs.uiuc.edu/pub/FWList/v1n13>>. Acesso em: 6 nov. 2006.
- JOHNSON, Ralph E. **How frameworks compare to other object-oriented reuse techniques**: Frameworks = (Components + Patterns). Communication of the ACM, New York-USA, v.40, n.10, p.39-42, Oct. 1997. Disponível em: <<http://portal.acm.org/citation.cfm?id=141943>>. Acesso em: 6 nov. 2006.
- JOHNSON, Ralph E.; FOOTE, Brian. **Designing Reusable Classes**. Journal of Object-Oriented Programming, New York-USA, v.1, n.2, p.22-35, Jun./Jul. 1988. Disponível em: <<http://www.laputan.org/drc/drc.html>>. Acesso em: 6 nov. 2006.
- KINOSHITA, Jorge; SALVADOR, Laís do Nascimento; MENEZES, Carlos Eduardo D. **CoGrOO – Um Corretor Gramatical para a língua portuguesa, acoplável ao OpenOffice**. In: Latin American Informatics Conference – CLEI2005, 31., 2005, Cali-Colômbia. **Anais...** Cali-Colômbia: Pontificia Universidad Javeriana-Cali, 2005. Disponível em: <[http://cogroo.incubadora.fapesp.br/portal/down/Doc/Artigo\\_Clei\\_2005.pdf](http://cogroo.incubadora.fapesp.br/portal/down/Doc/Artigo_Clei_2005.pdf)>. Acesso em: 18 nov. 2006.
- LEE, Richard C.; TEPFENHART, William M. **UML e C++: Guia Prático de Desenvolvimento Orientado a Objetos**. São Paulo: Makron Books, 2001. 550p.
- LEECH, Geoffrey. **The state of the art in corpus linguistics**: English Corpus Linguistics – Studies in Honour of Jan Svartvik. London-ENG: Longman, 1991.
- LEECH, Geoffrey; MACNERY, Anthony. **Corpus Annotation**: Linguistic Information from Computer Text Corpora. London-ENG: Longman, 1997. 281p.
- LINGUATECA. **WebJspell**: Questões e respostas relativas ao uso do Jspell. Disponível em: <<http://linguateca.di.uminho.pt/jspell/jsolhelp.pl>>. Acesso em: 15 nov. 2006.
- LUCENE. **Lucene!**. Disponível em: <<http://lucene.apache.org/index.html>>. Acesso em: 21 mai. 2007.
- MATTSSON, Michael. **Object-Oriented Frameworks**: A survey of methodological issues. Ronneby, Suécia: 1996. 122p. Tese (Doutorado) – Department of Computer Science and

Business Administration, University College of Karlskrona/Ronneby. Disponível em: <<http://citeseer.ist.psu.edu/article/mattsson96objectoriented.html>>. Acesso em: 7 nov. 2006.

MELO, Ana Cristina. **Desenvolvendo aplicações com UML 2.0: Do conceitual à implementação**. Rio de Janeiro: Brasport, 2004. 284p.

MENEZES, Carlos Eduardo D. de; NETO, João José. **Um Método Híbrido para a Construção de Etiquetadores Morfológicos, Aplicado à Língua Portuguesa, Baseado em Autômatos Adaptativos**. In: Conferência Ibero-americana em Sistemas, Cibernética e Informática – CISCI, 2002, Orlando-USA. **Anais...** Orlando-USA: CISCI, 2002, p.19-21. Disponível em: <[http://www.pcs.usp.br/~lta/artigos/menezes\\_cisci2002.pdf](http://www.pcs.usp.br/~lta/artigos/menezes_cisci2002.pdf)>. Acesso em: 15 nov. 2006.

NUNES, Maria da Graça V. et al. **A Construção de um Léxico para o Português do Brasil: Lições Aprendidas e Perspectivas**. In: Encontro para o processamento de português escrito e Falado, 2., 1996, Curitiba. **Anais...** Curitiba: CEFET-PR, 1996. p.61-70. Disponível em: <[www.icmc.sc.nilc.br/mdgvnune/download/curitilex.ps.gz](http://www.icmc.sc.nilc.br/mdgvnune/download/curitilex.ps.gz)>. Acesso em: 17 nov. 2006.

NUNES, Maria das Graças V. et al. **Introdução ao Processamento das Línguas Naturais**. Notas Didáticas do ICMC-USP, São Carlos, n.38, jun. 1999.

PARDO, Thiago Alexandre S.; RINO, Lucia Helena M.; NUNES, Maria das Graças V. **NeuralSumm: Uma abordagem Conexionista para a Sumarização Automática de Textos**. In: IV Encontro Nacional de Inteligência Artificial – ENIA. XXII Congresso Nacional da Sociedade Brasileira de Computação. **Anais...** Campinas: SBC, 2003. p.1-10. Disponível em: <<http://www.dc.ufscar.br/~lucia/articles/ENIA03-PardoEtAl.pdf>>. Acesso em: 14 nov. 2006.

PERCY, Carol E.; MEYER, Charles F.; LANCASHIRE, Ian. **Synchronic Corpus Linguistics: Papers from the sixteenth International Conference on English Language and Research on Computerized Corpora (ICAME 16)**. Amsterdam/Atlanta: Rodopi, 1996. 289p.

PRE, Wolfgang. **Design Patterns for Object-Oriented Software Development**. New York-USA: Addison-Wesley Publishing Company, 1995. 268p.

QUATRINI, Terry. **Modelagem Visual com Rational Rose 2000 e UML**. Rio de Janeiro: Editora Ciência Moderna Ltda, 2001. 206p.

RÉ, Reginaldo. **Um Processo para Construção de Frameworks a partir da Engenharia Reversa de Sistemas de Informação Baseados na Web: Aplicação ao Domínio de Leilões Virtuais**. São Paulo: 2002. 132p. Dissertação (Mestrado na Área de Ciências da Computação e Matemática Computacional) – Instituto de Ciências Matemáticas e de Computação, USP, 2002. Disponível em:

<<http://www.teses.usp.br/download.php/teses/disponiveis/55/55134/tde-20052003-120738/publico/dissert-f.pdf>>. Acesso em: 6 nov. 2006.

REED, Paul R. Jr. **Desenvolvendo Aplicativos com Visual Basic e UML**, São Paulo: Makron Books, 2000. 462p.

ROBERTS, Andrew. **JBootCat**. Disponível em: <<http://www.andyroberts.net/software/jbootcat/index.html>>. Acesso em: 15 nov. 2006.

ROCHE, Emmanuel; SCHABES, Yves. **Deterministic Part-of-Speech Tagging with Finite State Transducers**. Computational Linguistics, Cambridge-USA, v.21, n.2, p.227-253, jun. 1995. Disponível em: <[citeseer.ist.psu.edu/article/roche95deterministic.html](http://citeseer.ist.psu.edu/article/roche95deterministic.html)>. Acesso em: 15 nov. 2006.

- RUMBAUGH, James et al. **Modelagem e Projetos Baseados em Objetos**. Rio de Janeiro: Elsevier, 1994. 652p.
- RUSSEL, Stuart; NORVIG, Peter. **Inteligência Artificial**. Rio de Janeiro: Elsevier, 2004. 1040p.
- SANCHEZ, A. **Definicion e historia de los corpus**. In: A. SANCHEZ et al (org.). CUMBRE – Corpus Linguistico de Espanol Contemporaneo. Madrid-ESP: SGEL, 1995. p.7-24.
- SANT'ANNA, Victor Martins. **Cálculo de Referências Anafóricas Pronominais Demonstrativas na língua Portuguesa Escrita**. Porto Alegre: 2000. 100p. Dissertação (Mestrado em Ciência da Computação) - Faculdade de Informática, PUCRS, 2000. Disponível em: <[www.pucrs.br/uni/poa/info/pos/dissertacoes/arquivos/victor.pdf](http://www.pucrs.br/uni/poa/info/pos/dissertacoes/arquivos/victor.pdf)>. Acesso em: 17 nov. 2006.
- SARDINHA, Tony Berber. **Linguística de Corpus: Histórico e Problemática**. Delta, São Paulo, v.16, n.2, 2000. Disponível em: <[http://www.scielo.br/scielo.php?script=sci\\_arttext&pid=S0102-44502000000200005&lng=pt&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0102-44502000000200005&lng=pt&nrm=iso)>. Acesso em: 19 nov. 2006.
- SILVA, Luís Felipe S.; CHAPETTA, Wladimir Araújo. **Processamento de Linguagem Natural Aplicada à Inspeção de Documentos de Especificação de Requisitos de Software**. Disponível em: <[www.cos.ufrj.br/~ines/courses/cos740/leila/cos740/rel.pdf](http://www.cos.ufrj.br/~ines/courses/cos740/leila/cos740/rel.pdf)>. Acesso em: 17 nov. 2006.
- SILVA, Ricardo Pereira e. **Suporte ao desenvolvimento e Uso de frameworks e componentes**. Porto Alegre: 2000. 262p. Tese (Doutorado em Ciência da Computação) – Programa de Pós-Graduação em Computação, UFRGS, 2000. Disponível em: <<http://www.inf.ufsc.br/~ricardo/download/tese.pdf>>. Acesso em: 14 ago. 2006.
- SIMÕES, Alberto Manuel; ALMEIDA, José João. **Jspell.pm –Um módulo de Análise Morfológica para Uso em Processamento da Linguagem Natural**. In: Encontro da Associação Portuguesa de Linguística, 17., 2002, Lisboa-POR. Anais...Lisboa: 2002. p.485-495. Disponível em: <<http://alfarrabio.di.uminho.pt/~albie/publications/jspell.pm.pdf>>. Acesso em: 15 nov. 2006.
- SINCLAIR, John. **Corpus, Concordance, Collocation**. Oxford-USA: Oxford University Press, 1991. 197p.
- SNOWBALL. **Snowball**. Disponível em: <<http://www.snowball.tartarus.org/>>. Acesso em: 15 nov. 2006.
- TAGLIASSUCHI, Gustavo. **Mecanismo para Criação de Resumos: Google Summarizer**. Canoas: 2002. 40p. Monografia (Tecnólogo em Processamento de Dados) – Faculdade de Informática, ULBRA, 2002.
- TECHNORATI. **Technorati API Libraries**. Disponível em: <<http://www.technorati.com/developers/tools/libraries.html>>. Acesso em: 29 ago. 2006.
- VIEIRA, Renata; LIMA, Vera Lucua S. de. **Linguística computacional: Princípios e Aplicações**. In: Congresso da Sociedade Brasileira de Computação, 21., 2001, Fortaleza. Anais... Fortaleza: SBC, 2001, p.47-88. Disponível em: <[www.inf.unisinos.br/~renata/laboratorio/publicacoes/jaia12-vf.pdf](http://www.inf.unisinos.br/~renata/laboratorio/publicacoes/jaia12-vf.pdf)>. Acesso em: 15 nov. 2006.

VINHAES, Rêges Faria. **Estudo da Utilização de Técnica de Processamento de Linguagem Natural para Otimização de Tradutores Automáticos**. Rio Verde: 2005. 57p. Monografia (Bacharelado em Ciência da Computação) – Faculdade de Ciência da Computação, Universidade de Rio Verde, 2005.

YAHOO. **Yahoo! Search Web Services**. Disponível em: <<http://developer.yahoo.com/search/>>. Acesso em: 29 ago. 2006.

YOURDON, Edward; ARGILA, Carl. **Análise e Projeto Orientados a Objetos: Estudos de Caso**. São Paulo: Makron Books, 1999. 328p.