O Reconhecedor Gramatical Adaptativo Linguístico: Experimentos e Resultados Comparativos

A. T. Contier, D. Padovani e J. José Neto

Resumo— Este trabalho faz uma breve revisão dos conceitos de Tecnologia Adaptativa, apresentando seu mecanismo de funcionamento e seus principais campos de aplicação, destacando o forte potencial de sua utilização no processamento de linguagens naturais. Em seguida, são apresentados os conceitos de processamento de linguagem natural, ressaltando seu intricado comportamento estrutural. Por fim, é apresentado o Linguístico, um reconhecedor gramatical que utiliza autômatos e gramáticas adaptativas como tecnologia subjacente. O presente artigo é uma continuação de trabalhos anteriores, com uma nova seção para apresentar os experimentos realizados e os resultados obtidos comparados aos do estado da arte.

Palavras Chave — Autômatos Adaptativos, Processamento de Linguagem Natural, Reconhecedores Gramaticais, Gramáticas Livres de Contexto, Gramáticas Adaptativas.

I. AUTÔMATOS ADAPTATIVOS

AUTÔMATO adaptativo é uma máquina de estados à qual são impostas sucessivas alterações resultantes da aplicação de ações adaptativas associadas às regras de transições executadas pelo autômato [1]. Dessa maneira, estados e transições podem ser eliminados ou incorporados ao autômato em decorrência de cada um dos passos executados durante a análise da entrada. De maneira geral, pode-se dizer que o autômato adaptativo é formado por um dispositivo convencional, não adaptativo, e um conjunto de mecanismos adaptativos responsáveis pela auto modificação do sistema.

O dispositivo convencional pode ser uma gramática, um autômato, ou qualquer outro dispositivo que respeite um conjunto finito de regras estáticas. Este dispositivo possui uma coleção de regras, usualmente na forma de cláusulas if-then, que testam a situação corrente em relação a uma configuração específica e levam o dispositivo à sua próxima situação. Se nenhuma regra é aplicável, uma condição de erro é reportada e a operação do dispositivo, descontinuada. Se houver uma única regra aplicável à situação corrente, a próxima situação do dispositivo é determinada pela regra em questão. Se houver mais de uma regra aderente à situação corrente do dispositivo, as diversas possíveis situações seguintes são tratadas em paralelo e o dispositivo exibirá uma operação não determinística. Os mecanismos adaptativos são formados por três tipos de ações adaptativas elementares: consulta (inspeção do conjunto de regras que define o dispositivo), exclusão (remoção de alguma regra) e inclusão (adição de uma nova

Autômatos adaptativos apresentam forte potencial de aplicação ao processamento de linguagens naturais, devido à facilidade com que permitem representar fenômenos linguísticos complexos tais como dependências de contexto.

Adicionalmente, podem ser implementados como um formalismo de reconhecimento, o que permite seu uso no préprocessamento de textos para diversos usos, tais como: análise sintática, verificação de sintaxe, processamento para traduções automáticas, interpretação de texto, corretores gramaticais e base para construção de sistemas de busca semântica e de aprendizado de línguas auxiliados por computador.

Diversos trabalhos confirmam a viabilidade prática da utilização de autômatos adaptativos para processamento da linguagem natural. É o caso, por exemplo, de [2], que mostra a utilização de autômatos adaptativos na fase de análise sintática; [3] que apresenta um método de construção de um analisador morfológico e [4], que apresenta uma proposta de autômato adaptativo para reconhecimento de anáforas pronominais segundo algoritmo de Mitkov.

II. PROCESSAMENTO DA LINGUAGEM NATURAL: REVISÃO DA LITERATURA

O processamento da linguagem natural requer o desenvolvimento de programas que sejam capazes de determinar e interpretar a estrutura das sentenças em muitos níveis de detalhe. As linguagens naturais exibem um intricado comportamento estrutural visto que são profusos os casos particulares a serem considerados. Uma vez que as linguagens naturais nunca são formalmente projetadas, suas regras sintáticas não são simples nem tampouco óbvias e tornam, portanto, complexo o seu processamento computacional. Muitos métodos são empregados em sistemas processamento de linguagem natural, adotando diferentes paradigmas, tais como métodos exatos, aproximados, prédefinidos ou interativos, inteligentes ou algorítmicos [5]. Independentemente do método utilizado, o processamento da linguagem natural envolve as operações de análise léxicomorfológica, análise sintática, análise semântica e análise pragmática [6]. A análise léxico-morfológica procura atribuir uma classificação morfológica a cada palavra da sentença, a partir das informações armazenadas no léxico [7]. O léxico ou dicionário é a estrutura de dados contendo os itens lexicais e as informações correspondentes a estes itens. Entre as informações associadas aos itens lexicais, encontram-se a categoria gramatical do item, tais como substantivo, verbo e adjetivo, e os valores morfossintático-semânticos, tais como gênero, número, grau, pessoa, tempo, modo, regência verbal ou nominal. Na etapa de análise sintática, o analisador verifica se uma sequência de palavras constitui uma frase válida da língua, reconhecendo-a ou não. O analisador sintático faz uso de um léxico e de uma gramática, que define as regras de combinação dos itens na formação das frases. Nos casos nos quais há a necessidade de interpretar o significado de um texto, a análise léxico-morfológica e a análise sintática não são suficientes, sendo necessário realizar um novo tipo de operação, denominada análise semântica [7]. Na análise semântica procura-se mapear a estrutura sintática para o domínio da aplicação, fazendo com que a estrutura ganhe um significado. O mapeamento é feito identificando as propriedades semânticas do léxico e o relacionamento semântico entre os itens que o compõe [8]. Já a análise pragmática procura reinterpretar a estrutura que representa o que foi dito para determinar o que realmente se quis dizer. Inserem-se nessa categoria as relações anafóricas, correferências, determinações, focos ou temas, dêiticos e elipses [9].

Em [10] são apresentados trabalhos de pesquisas em processamento de linguagem natural para a Língua Portuguesa tais como o desenvolvido pelo Núcleo Interinstitucional de Linguística Aplicada (NILC) no desenvolvimento de ferramentas para processamento de linguagem natural; o projeto VISL - Visual Interactive Syntax Learning, sediado na Universidade do Sul da Dinamarca, que engloba o desenvolvimento de analisadores morfossintáticos para diversas línguas, entre as quais o português; e o trabalho de resolução de anáforas desenvolvido pela Universidade de Santa Catarina. A tecnologia adaptativa também tem contribuído com trabalhos em processamento da linguagem natural. Em [11], são apresentadas algumas das pesquisas desenvolvidas pelo Laboratório de Linguagens e Tecnologia Adaptativa da Escola Politécnica da Universidade de São Paulo: um etiquetador morfológico, um estudo sobre processos de análise sintática, modelos para tratamento de não determinismos e ambiguidades, e um tradutor texto voz baseado em autômatos adaptativos.

III. RECONHECEDOR ADAPTATIVO: SUPORTE TEÓRICO LINGUÍSTICO

A Moderna Gramática Brasileira de Celso Luft [12] foi escolhida como suporte teórico linguístico do reconhecedor aqui proposto. A escolha foi feita em função da forma clara e precisa com que Luft categoriza os diversos tipos de sentenças de língua portuguesa, diferenciando-se das demais gramáticas que priorizam a descrição da língua em detrimento da análise estrutural da mesma. Luft diz que a oração é moldada por padrões frasais ou oracionais, compostos por elementos denominados sintagmas (Tabela 1).

TABELA 1. ELEMENTOS FORMADORES DE SINTAGMAS

Sintagmas				
Substantivo	Quantitativos+Pronomes Adjetivos+ Sintagma Adjetivo1+Substantivo+ Sintagma Adjetivo2+ Sintagma Preposicional+ Oração Adjetiva			
Verbal	Pré-verbais+ Verbo Auxiliar+ Verbo Principal			

TABELA 1. CONTINUAÇÃO

Sintagmas			
Adjetivo	Advérbio de Intensidade+		
	Adjetivo+		
	Sintagma Preposicional		
Adverbial	Advérbio de Intensidade+		
	Adverbio+		
	Sintagma Preposicional		
Preposicional	Preposição+		
	Sintagma Substantivo		

Sintagma é qualquer constituinte imediato da oração, podendo exercer papel de sujeito, complemento (objeto direto e indireto), predicativo e adjunto adverbial. É composto por uma ou mais palavras, sendo que uma é classificada como núcleo e as demais como dependentes. As palavras dependentes podem estar localizadas à esquerda ou à direita do núcleo. Luft utiliza os seguintes nomes e abreviaturas:

- 1. Sintagma substantivo (SS): núcleo é um substantivo;
- 2. Sintagma verbal (SV): núcleo é um verbo;
- 3. Sintagma adjetivo (Sadj): núcleo é um adjetivo;
- 4. Sintagma adverbial (Sadv): núcleo é um advérbio;
- 5. Sintagma preposicional (SP): é formado por uma preposição (Prep) mais um SS.
- 6. Vlig: verbo de ligação
- 7. Vi: verbo intransitivo
- 8. Vtd: verbo transitivo direto
- 9. Vti: verbo transitivo indireto
- 10. Vtdi: verbo transitivo direto e indireto 11. Vt-pred: verbo transitivo predicativo

Um padrão oracional é determinado pelos tipos de sintagmas e pela sequência em que aparecem. Por exemplo, o padrão oracional SS Vlig SS, indica que a frase é composta por um sintagma substantivo, seguido de um verbo de ligação e de outro sintagma substantivo. Os padrões são classificados em cinco tipos (Tabela 2):

- 1. Padrões pessoais nominais: Neste caso, existe sujeito e o núcleo do predicado é um nome (substantivo, adjetivo, advérbio) ou um pronome (substantivo, adjetivo, advérbio). O verbo, nesses casos, é chamado de verbo de ligação (Vlig).
- 2. Padrões pessoais verbais: São aqueles nos quais existe o sujeito e o núcleo do predicado é um verbo. O verbo pode ser transitivo direto (Vtd), transitivo indireto (Vti), transitivo direto e indireto (Vtdi), e intransitivo (Vi). Se o verbo for transitivo direto (Vtd), o complemento será um objeto direto; se o verbo for transitivo indireto (Vti), o complemento será um objeto indireto; se o verbo for transitivo direto e indireto (Vtdi), o complemento será um objeto direto e um indireto; se o verbo for intransitivo (Vi), não há complemento.
- Padrões Pessoais Verbo-Nominais: Neste caso, existe o sujeito e o núcleo do predicado é um verbo transitivo predicativo (Vt-pred), cujo complemento é um objeto direto e um predicativo do objeto.

- Padrões Impessoais Nominais: Ocorrem quando não existe sujeito e o núcleo do predicado é um nome (substantivo, adjetivo, advérbio) ou um pronome (substantivo, adjetivo, advérbio).
- Padrões Impessoais Verbais: Neste caso, não existe sujeito e o núcleo do predicado é um verbo.

TABELA 2 PADRÕES ORACIONAIS DE LUFT

Padrões Pessoais Nominais						
SS	Vlig	SS				
SS	Vlig	Sadj				
SS	Vlig	Sadv				
SS	Vlig	SP				
Padrões P	essoais Verb	ais	l.	•		
SS	Vtd	SS				
SS	Vti	SP				
SS	Vti	Sadv				
SS	Vti	SP	SP			
SS	Vtdi	SS	SP			
SS	Vtdi	SS	Sadv			
SS	Vtdi	SS	SP	SP		
SS	Vi					
Padrões P	essoais Verbe	o-Nominais		•		
SS	Vtpred	SS	SS			
SS	Vtpred	SS	Sadj			
SS	Vtpred	SS	SP			
SS	Vtpred	SS	Sadv			
SS	Vtpred	SS				
SS	Vtpred	Sadj				
SS	Vtpred	SP				
Padrões In	npessoais No					
	Vlig	SS				
	Vlig	Sadj				
	Vlig	Sadv				
D 1 2 7	Vlig SP					
Padrões Impessoais Verbais						
	Vtd	SS				
	Vti	SP				
	Vi					

Os sintagmas e os padrões oracionais de Luft foram mapeados em uma gramática para que pudessem ser processados computacionalmente, ficando da seguinte forma:

 $S \rightarrow [Conec][SS]SV[Conec]$

Conec \rightarrow S

 $SS \rightarrow [Sadj] SS [Sadj | SP]$

 $SS \rightarrow [Quant | PrA] (Sc | Sp | PrPes)$

SV→ [Neg] [Aux | PreV] (Vlig | Vtd| Vti | Vtdi| Vi)

[SS | Sadj | Sadv | SP] [SS | Sadj | Sadv | SP] [SP]

 $SP \rightarrow Prep (SS | Sadj)$

Sadj → Sadj [SP]

Sadj \rightarrow [Adv] Adj

Sadv → Sadv [SP]

Sadv → [Adv] Adv

 $PrA \boldsymbol{\rightarrow} \quad Ind \mid ArtDef \mid ArtInd \mid Dem \mid Pos$

Sendo:

S – Sentença

SS – Sintagma substantivo

SV – Sintagma verbal

SP – Sintagma preposicional

SN – Sintagma nominal

Sadv – Sintagma adverbial

Sadj - Sintagma adjetivo

Adv - adverbio

Adj – adjetivo

ArtDef - artigo definido

ArtInd - artigo indefinido

Aux – Partícula auxiliar (apassivadora ou pré-verbal)

Conec – Conector (conjunção ou pronome relativo)

Dem – pronome demonstrativo indefinido

Ind – pronome indefinido

Neg – partícula (negação)

PrA – pronome adjetivo

PrPes – pronome pessoal

Prep - preposição

Ouant – numeral

Sc – substantivo comum

Sp – substantivo próprio

V- verbo

Vlig - verbo de ligação

Vi – verbo intransitivo

Vtd - verbo transitivo direto

Vti - verbo transitivo indireto

Vtdi – verbo transitivo direto e indireto

IV. PROPOSTA DE UM RECONHECEDOR GRAMATICAL

O Linguístico é uma proposta de reconhecedor gramatical composto de cinco módulos sequenciais que realizam cada qual um processamento especializado, enviando o resultado obtido para o módulo seguinte, tal como ocorre em uma linha de produção, até que o texto esteja completamente analisado (Fig.1).

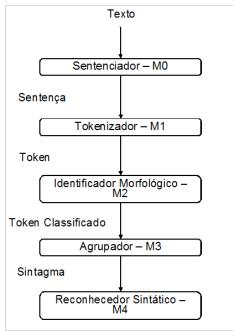


Figura 1. Estrutura do Linguístico.

O primeiro módulo, denominado Sentenciador, recebe um texto e realiza um pré-processamento, identificando os caracteres que possam indicar final de sentença, palavras abreviadas e palavras compostas, e eliminando aspas simples e duplas. Ao final, o Sentenciador divide o texto em supostas sentenças, para análise individual nas etapas seguintes.

O segundo módulo, denominado *Token*izador, recebe as sentenças identificadas na etapa anterior e as divide em *tokens*, considerando, neste processo, abreviaturas, valores monetários, horas e minutos, numerais arábicos e romanos, palavras compostas, nomes próprios, caracteres especiais e de pontuação final. Os *tokens* são armazenados em estruturas de dados (*arrays*) e enviados um a um para análise do módulo seguinte.

O terceiro módulo, chamado Identificador Morfológico, foi concebido para utilizar tecnologias adaptativas (Fig.2.1). O Identificador Morfológico é composto por um Autômato Mestre e um conjunto de máquinas especialistas que acessam bases de dados de regras de acesso às classificações morfológicas. A prioridade é obter as classificações do corpus Bosque [13]; caso o termo procurado não seja encontrado, o Identificador Morfológico procura por substantivos, adjetivos e verbos, no formato finito e infinito (flexionados e não flexionados), através de uma máquina de formação e identificação de palavras, que usa, como base, o vocabulário do TeP2.0 [14]; por fim, o Identificador Morfológico procura por termos invariáveis, ou seja, termos cuja classificação morfológica é considerada estável pelos linguistas, tais como, conjunções, preposições e pronomes, no caso, extraídos no léxico do Portal São Francisco [15].

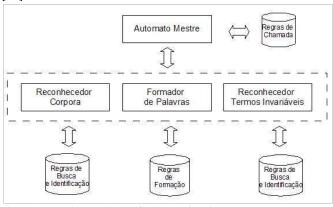


Figura 2.1. Arquitetura do Identificador Morfológico.

O Autômato Mestre (Fig.2.2) é responsável pelo sequenciamento das chamadas às máquinas, de acordo com um conjunto de regras cadastradas em base de dados. Ele inicia o processamento recebendo o *token* da coleção de *tokens* criada pela classe Texto. Em seguida, antes de processá-lo, o autômato se modifica através de uma função adaptativa, criando uma máquina de processamento (M1), uma transição entre o estado 1 e M1, e uma transição que aguarda o estado final da máquina M1. O *token* é passado para M1 e armazenado em uma pilha. Quando M1 chega ao estado final, o autômato se modifica novamente, criando uma nova máquina M2, o estado 2 e as transições correspondentes. Caso o estado final de M1 seja de aceitação, o processo é finalizado no estado 2, caso contrário a máquina M2 é chamada, passando o *token* armazenado na pilha.

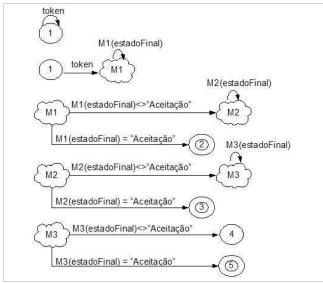


Figura 2.2. Estrutura Adaptativa do Autômato Mestre.

O processo se repete quando M2 chega ao final do processamento, com a criação da máquina M3, do estado 3 e das transições correspondentes. Um novo ciclo se repete, e se o estado final de M3 é de aceitação, o autômato transiciona para o estado 5, de aceitação, caso contrário, ele vai para o estado 4 de não aceitação. M1, M2 e M3 representam, respectivamente, as máquinas do Reconhecedor de Corpus, do Formador de Palavras e do Reconhecedor de Termos Invariáveis. Portanto, caso o Autômato Mestre encontre a classificação morfológica ao final de M1, ele não chama M2; caso encontre em M2, não chama M3 e, caso também não encontre em M3, ele informa aos demais módulos do Linguístico que não há classificação morfológica para o termo analisado. Embora o Autômato Mestre processe um número pré-determinado de transições de estado e inicialmente não apresente variação do fluxo, optou-se pelo uso da tecnologia adaptativa devido à flexibilidade com que regras podem ser incluídas ou alteradas, o que é particularmente útil, caso seja necessário usar novas técnicas de reconhecimento.

As máquinas M1, M2 e M3 também foram projetadas de acordo com a tecnologia adaptativa. A Máquina M1 usa um autômato adaptativo que se auto modifica de acordo com o tipo de *token* que está sendo analisado: palavras simples, palavras compostas, números, valores e símbolos. A Fig. 2.3 apresenta a estrutura adaptativa de M1.

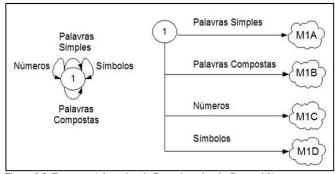


Figura 2.3. Estrutura Adaptativa do Reconhecedor de Corpus M1.

Inicialmente o autômato é composto por um único estado e por transições para ele mesmo (Fig.2.3, à esquerda). Ao identificar o tipo de *token* que será analisado (obtido no processo de *token*ização), o autômato cria máquinas e as transições correspondentes. As alternativas de configuração são apresentadas na Fig.2.3, à direita. As submáquinas M1A, M1B, M1C e M1D reconhecem os *tokens* através de outro tipo de autômato que processa as unidades conforme apresentado na Fig. 2.4.

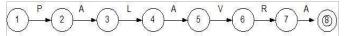


Figura 2.4. Reconhecedor de Palavras Simples.

No exemplo apresentado na Fig.2.4, o *token* "Palavra" é processado pela submáquina M1A; se o processamento terminar em um estado de aceitação, o *token* é reconhecido. As submáquinas M1A, M1B, M1C e M1D são criadas previamente por um programa que lê o Corpus e o converte em autômatos finitos determinísticos. A estrutura de identificação morfológica é composta por um par [chave, valor], no qual a chave é o estado de aceitação do elemento lexical e, o valor, sua classificação morfológica. No exemplo apresentado, a chave do item lexical "palavra" seria o estado "8", e, o valor, a classificação morfológica, H+n (substantivo, núcleo de sintagma nominal), proveniente do Corpus Bosque.

O Formador de Palavras, máquina M2, também é montado previamente por um programa construtor, usando o vocabulário do TeP2.0 como léxico de formas previamente construídas e o conjunto de regras de prefixação, sufixação e regressão verbal descrito por Margarida Basílio [16] para construir novas formas. As regras de formação são mapeadas previamente ao conjunto de palavras ao qual podem ser aplicadas, evitando repetições e permitindo a seleção das formas pertinentes no momento em que a palavra em questão é analisada. No entanto, são necessários alguns cuidados para evitar a criação de estruturas que aceitem palavras inexistentes. A Fig. 2.5 apresenta um exemplo de autômato no qual são aplicados os prefixos "a" e "per", e os sufixos "ecer" e "ar" (derivação parassintética) ao radical "noit" do substantivo noite, que faz parte do TeP2.0.

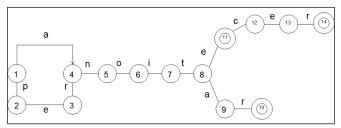


Figura 2.5. Autômato Formador de Palavras.

No exemplo apresentado, as palavras "anoitecer", "pernoitar" e "anoitar" (sinônimo de "anoitecer") existem no léxico do português do Brasil. Já pernoitecer é uma combinação que não existe na língua portuguesa. Margarida Basílio diz que algumas combinações não são aceitas simplesmente porque já existem outras construções consagradas pelo uso [17]. Para reduzir o risco de aceitar derivações inexistentes, o processo de construção do autômato

restringe as possíveis formações, utilizando apenas as regras que a autora destaca como sendo mais prováveis. É o caso de nominalização de verbos com o uso dos sufixos –ção , –mento e – da. Em [16], Basílio cita que o sufixo –ção é responsável por 60% das formações regulares, enquanto o sufixo –mento é responsável por 20% destas formações. Já o sufixo –da é, via de regra, usado em nominalizações de verbos de movimento, tais como, entrada, saída, partida, vinda, etc.

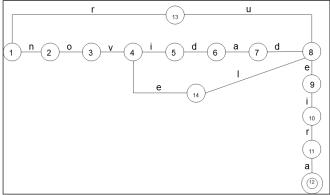


Figura 2.6. Autômato Formador de Palavras.

Outra característica do construtor do autômato é representar os prefixos e sufixos sempre pelo mesmo conjunto de estados e transições, evitando repetições que acarretariam o consumo desnecessário de recursos computacionais. A Fig. 2.6 apresenta exemplo de palavras formadas por reutilização de estados e transições na derivação das palavras rueira, novidadeira e noveleira. Os estados e transições usados para representar o sufixo "eira", usado para designar são os mesmos nas três derivações.

A máquina M2 também reconhece palavras flexionadas, obtidas, no caso de substantivos e adjetivos, através da aplicação de sufixos indicativos de gênero, número e grau aos radicais do vocabulário TeP2.0. A Fig. 2.7 apresenta um exemplo de autômato usado para a formação das formas flexionadas do substantivo menino. Foram adicionados ao radical "menin" as flexões "o(s)", "a(s)", "ão", "ões", "onona(s)", "inho(s)" e"inha(s)".

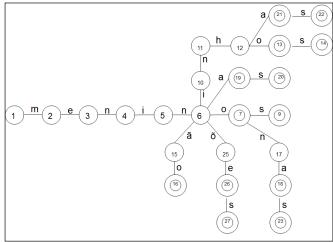


Figura 2.7. Autômato Formador de Flexões Nominais.

No caso de verbos, foi criada uma estrutura de estados e transições para representar as flexões de tempo, modo, voz e pessoa, obtendo-se, assim, as respectivas conjugações. A Fig. 2.8 apresenta um exemplo de autômato usado para a formação das formas flexionadas do presente do indicativo do verbo andar. Foram adicionados ao radical "and" as flexões "o", "as", "a", "andamos", "ais" e "andam".

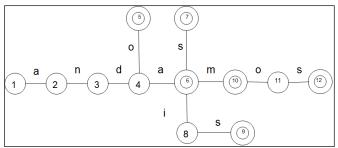


Figura 2.8. Autômato Formador de Flexões Verbais.

A estrutura de armazenamento da máquina M2 também é composta por um par [chave, valor], no qual a chave é o estado de aceitação do elemento lexical e, o valor, sua classificação morfológica, acrescida da origem do termo, indicando que ele foi gerado pelo construtor. Por exemplo, o termo "novidadeira" seria classificado como H+n+C – substantivo, núcleo de sintagma nominal, gerado pelo construtor.

Já a máquina M3 é um autômato que varia em função do tipo de termo (conjunções, preposições e pronomes) e utiliza uma estrutura arbórea similar a submáquina M1A. A Fig. 2.9 apresenta um exemplo de autômato usado no reconhecimento de termos deste domínio.



Figura 2.9. Autômato Reconhecedor de Conjunções, Preposições e Pronomes.

A estrutura de armazenamento da máquina M3 é composta por um par [chave, valor], no qual a chave é o estado de aceitação do elemento lexical e, o valor, sua classificação morfológica, acrescida da origem do termo, indicando que ele faz parte da base de dados de apoio do Portal São Francisco. Por exemplo, o termo "embora" seria classificado como cj+Ba –conjunção da base de dados de apoio.

O Autômato Mestre também é responsável por desambiguar as classificações morfológicas e informar ao Agrupador apenas as que forem mais adequadas. Como as palavras podem ter mais do que uma classificação morfológica, é necessário utilizar uma técnica para selecionar aquela que é mais apropriada para o contexto analisado. Por exemplo, a palavra "casa" pode ser classificada como substantivo comum, feminino, singular ou como verbo flexionado na 3ª pessoa do singular no tempo presente, modo indicativo. No entanto, tendo em vista o contexto em que palavra se encontra, é possível escolher a classificação mais provável. Por exemplo, se a palavra "casa" vier precedida de um artigo definido, é mais provável que ela seja um substantivo; já se a palavra antecessora for um substantivo próprio, é mais provável que "casa" seja um verbo.

Collins [18] apresenta um modelo estatístico que leva em consideração 2 classificações anteriores à palavra analisada para

fazer a desambiguação, criando distribuições nas quais a etiqueta de máxima probabilidade é o resultado da função:

 $P = \max p(E,S)$, sendo:

E = etiquetas

S = palavras da sentença

p = probabilidade de E, dado S

max = máxima probabilidade

Por exemplo, para etiquetar a frase "O advogado entrevista a testemunha", a função receberia as palavras da sentença e as sequências de etiquetas possíveis para elas (obtidas a partir de um Corpus de testes), calculando, então, a sequência de maior probabilidade. No entanto, Collins alerta que a complexidade para executar tal função inviabiliza sua utilização, pois ela cresce exponencialmente em função do número de palavras da sentença (Em uma sentença de "n" palavras e "K" possíveis classificações, existiriam |K|ⁿ possíveis etiquetas de sequencia).

Para evitar o problema da complexidade da execução, Collins propõe o uso do algoritmo Viterbi [19]. O algoritmo Viterbi é usado para encontrar a sequência mais provável de estados ocultos que resultam da observação de uma sequência de eventos. No caso da identificação das etiquetas morfológicas, os eventos são as palavras do texto analisado e os estados são as etiquetas de identificação morfológica. Para uma dada sentença ou sequencia de palavras, o algoritmo Viterbi escolhe a etiqueta que maximiza a seguinte fórmula:

P(token|etiqueta) * P(etiqueta| n etiquetas anteriores) (1)

O algoritmo Viterbi pode ser implementado por meio de um autômato finito. No entanto, um autômato com estas características poderia ficar muito grande e, consequentemente, consumir elevados recursos computacionais, devido ao tamanho do léxico usado para montá-lo. Uma alternativa para evitar este tipo de problema, é usar a tecnologia adaptativa. A Fig. 2.10 apresenta um exemplo da estrutura do autômato que implementa o algoritmo Viterbi usado pelo desambiguador morfológico do Linguístico. A sequência de tokens analisados é "A casa". O passo apresentado mostra como o autômato seleciona a etiqueta do token "casa". A primeira etiqueta foi identificada como "Início" e a segunda como "Artigo". Supõese, no exemplo, que além de "casa", "bicicleta" e "cidade" sejam tokens possíveis, dadas as sequências de etiquetas analisadas (Início-> Artigo->Substantivo e Início -> Artigo -> Verbo).

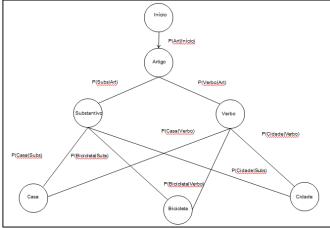


Figura 2.10. Autômato Adaptativo Viterbi.

O autômato monta dinamicamente os estados e transições de acordo com as etiquetas e probabilidades, identificando, ao final do processamento, a etiqueta mais provável. No exemplo apresentado, a seleção é feita através da comparação das seguintes sequências:

```
Início -> Artigo -> Substantivo -> Casa
Início -> Artigo -> Substantivo -> Bicicleta
Início -> Artigo -> Substantivo -> Cidade
Início -> Artigo -> Verbo -> Casa
Início -> Artigo -> Verbo -> Bicicleta
Início -> Artigo -> Verbo -> Cidade
```

Ao término das comparações, o autômato seleciona a etiqueta que maximiza a probabilidade de ocorrência de acordo com a fórmula geral de Viterbi (1), no caso, "Substantivo".

O quarto módulo, denominado Agrupador é composto de um autômato, responsável pela montagem dos sintagmas a partir de símbolos terminais da gramática e um bigrama, responsável pela montagem dos sintagmas a partir de não-terminais (Fig.3). Inicialmente, o Agrupador recebe do Identificador as classificações morfológicas dos *tokens* e as agrupa em sintagmas de acordo com as regras de Luft. Neste processo são identificados sintagmas nominais, verbais, preposicionais, adjetivos e adverbiais Para isso, o Agrupador utiliza um autômato adaptativo cuja configuração completa é definida da seguinte forma:

```
Estados = \{1, 2, 3, 4, SS, SP, V, Sadj, Sadv, A\}, Onde:
```

1,2,3 e 4 = Estados Intermediários

SS, SP, V Sadj, Sadv = Estados nos quais houve formação de sintagmas, sendo:

SS= Sintagma substantivo

SP = Sintagma preposicional

V = Verbo ou locução verbal

Sadj = Sintagma adjetivo

Sadv = Sintagma adverbial

A = Estado após o processamento de um ponto final

Tokens = {art, num, n, v, prp, pron, conj, adj, adv, rel, pFinal, sClass}, onde:

art = artigo, num = numeral

n = substantivo, v = verbo

prp = preposição, pron = pronome

conj = conjunção, adj = adjetivo

adv = advérbio, rel = pronome relativo

pFinal = ponto final, sClass = sem classificação

Estados de Aceitação = {SS, SP, V, Sadj, Sadv, A}

Estado Inicial = $\{1\}$

Função de Transição = {(Estado, *Token*)→Estado}, sendo:

 $\{(1, \operatorname{art}) \rightarrow 2, (2, \operatorname{art}) \rightarrow 2, (3, \operatorname{art}) \rightarrow 3\}$

 $(1, \text{num}) \rightarrow \text{Sadv}, (2, \text{num}) \rightarrow 2, (3, \text{num}) \rightarrow 3$

 $(1, n) \rightarrow SS, (2, n) \rightarrow SS, (3, n) \rightarrow SP$

 $(1, v) \rightarrow SV, (2, v) \rightarrow SV, (3, v) \rightarrow SP$

 $(1, prp) \rightarrow 3, (2, prp) \rightarrow 2, (3, prp) \rightarrow 3$

 $(1, prop) \rightarrow SS, (2, prop) \rightarrow SS, (3, prop) \rightarrow SP$

 $(1, pron) \rightarrow SS, (2, pron) \rightarrow SS, (3, pron) \rightarrow SP$

 $(1, \operatorname{conj}) \rightarrow \operatorname{conj}, (2, \operatorname{conj}) \rightarrow \emptyset, (3, \operatorname{conj}) \rightarrow \emptyset$

 $(1, adj) \rightarrow Sadj, (2, adj) \rightarrow Sadj, (3, adj) \rightarrow 3$

 $(1, adv) \rightarrow Sadv, (2, adv) \rightarrow 2, (3, adv) \rightarrow 3$

 $(1, \text{rel}) \rightarrow \text{conj}, (2, \text{rel}) \rightarrow \emptyset, (3, \text{rel}) \rightarrow \text{conj}$ $(1, \text{pFinal}) \rightarrow A, (2, \text{pFinal}) \rightarrow \emptyset, (3, \text{pFinal}) \rightarrow \emptyset$

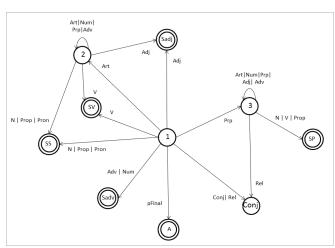


Figura 3. Configuração Completa do Autômato Construtor de Sintagmas.

Por exemplo, no caso de sintagmas substantivos teríamos: SS → [Quant | PrA] (Sc | Sp | PrPes)

Pela regra acima, o conjunto de *tokens* "A" e "casa" formam um sintagma substantivo, da seguinte forma:

PrA = "A" (artigo definido)

Sc = "casa" (substantivo comum)

Da direita para esquerda, são realizadas as seguintes transições:

PrA Sc \rightarrow SS; A Sc \rightarrow SS; A casa \rightarrow SS

Já o Agrupador recebe o *token* "A", identificado pelo *Token*izador como artigo definido, e se movimenta do estado 1 para o estado 2. Ao receber o *token* "casa", identificado como substantivo comum, ele se movimenta do estado 2 para o estado SS, que é um estado de aceitação. Neste momento o Agrupador armazena a cadeia "A casa" e o símbolo "SS" em uma pilha e reinicializa o autômato preparando-o para um novo reconhecimento. Em um passo seguinte, o Agrupador usa o bigrama para comparar um novo sintagma com o último sintagma formado, visando identificar elementos mais altos na hierarquia. Para isso ele usa a matriz apresentada na Tabela 3. A primeira coluna da matriz indica o último sintagma formado (US) e a primeira linha, o sintagma atual (SA). A célula resultante apresenta o novo nó na hierarquia.

TABELA 3.
MATRIZ DE AGRUPAMENTO DE SINTAGMAS

SA US	SS	SP	V	Sadv	Sadj	Conj
SS	SS	SS	-	-	SS	-
SP	SP	-	-	-	-	-
V	-	-	V	-	-	-
Sadv	-	-	-	Sadv	Sadj	-
Sadj	SS	Sadj	1	-	Sadj	-
Conj	-	-	-	-	-	Conj

Esta técnica foi usada para tratar as regras gramaticais nas quais um sintagma é gerado a partir da combinação de outros,

como é o caso da regra de formação de sintagmas substantivos: SS \rightarrow [Sadj | SP]. Por esta regra, os sintagmas substantivos são formados por outros sintagmas substantivos precedidos de um sintagma adjetivo e seguidos de um sintagma adjetivo ou um sintagma preposicional. No exemplo anterior, supondo que os próximos 2 tokens fossem "de" e "madeira", após a passagem pelo autômato, o Agrupador formaria um sintagma SP. Considerando que na pilha ele tinha armazenado um SS, após a passagem pelo bigrama, e de acordo com a Tabela 3, o sintagma resultante seria um SS e o conteúdo que o compõe seria a combinação dos textos de cada sintagma que o originou. Caso não haja agrupamentos possíveis, o Agrupador envia o último sintagma formado para análise do Reconhecedor Sintático e movimenta o sintagma atual para a posição de último sintagma no bigrama, repetindo o processo com o próximo sintagma.

O quinto e último módulo, denominado Reconhecedor Sintático, recebe os sintagmas do módulo anterior e verifica se estão sintaticamente corretos de acordo com padrões oracionais de Luft. O Reconhecedor Sintático utiliza um autômato adaptativo que faz chamadas recursivas sempre que recebe conjunções ou pronomes relativos, armazenando, em uma estrutura de pilha, o estado e a cadeia de sintagmas reconhecidos até o momento da chamada. Caso o Reconhecedor Sintático não consiga se movimentar a partir do sintagma recebido, ele gera um erro e retorna o ponteiro para o último sintagma reconhecido, finalizando a instância do autômato recursivo e retornando o processamento para aquela que a inicializou. Esta, por sua vez, retoma posição em que se encontrava antes da chamada e continua o processamento até o final da sentença ou até encontrar uma nova conjunção, situação na qual o processo se repete.

A configuração completa do autômato é definida da seguinte forma:

Estados = { 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27}

Tokens = {SS, SP, Vli, Vi, Vtd, Vti, Vtdi, Sadj, Sadv, Conj, A}

Estados de Aceitação = { 4, 5, 6, 9, 12, 13,14, 15, 17, 18, 19, 21, 22, 24, 25, 26, 27}

Estado Inicial = $\{1\}$

Função de Transição = {(Estado, *Token*)→Estado}, sendo:

 $\{(1, SS) \rightarrow 2, (2, Vti) \rightarrow 3, (3, SP) \rightarrow 4, (4, SP) \rightarrow 4,$

- $(3, Sadv) \rightarrow 5, (2, Vi) \rightarrow 6, (2, Vtdi) \rightarrow 7$
- $(7, SS) \rightarrow 8, (8, SP) \rightarrow 9, (9, SP) \rightarrow 9, (8, Sadv) \rightarrow 10,$
- $(2, Vlig) \rightarrow 11, (11, SP) \rightarrow 12, (11, Sadv) \rightarrow 13,$
- $(11, Sadj) \rightarrow 14, (11, SS) \rightarrow 15, (2, Vtd) \rightarrow 16, (16, SS) \rightarrow 17,$
- $(2, Vtpred) \rightarrow 18, (18, SP) \rightarrow 19, (18, Sadj) \rightarrow 20$
- $(18, SS) \rightarrow 21, (21, SS) \rightarrow 22, (21, Sadj) \rightarrow 23, (21, Sadv) \rightarrow 24, (21, SP) \rightarrow 25$

Pilha = {[Texto, Sintagma, Estado]}

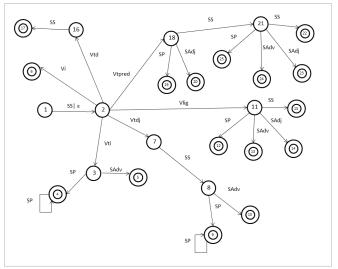


Figura 4.1. Configuração Completa do Reconhecedor Sintático.

No entanto, para que a análise sintática seja feita, não são necessárias todas as ramificações da configuração completa do autômato (Fig. 4.1). Por exemplo, quando se transita um verbo de ligação a partir do estado 2, o autômato vai para o estado 11 e todas as demais ramificações que partem deste estado para os estados 3, 7, 16 e 18, não são usadas. Com a tecnologia adaptativa, é possível criar dinamicamente os estados e transições do autômato em função dos tipos de verbos, evitando manter ramificações que não são usadas.

A Fig. 4.2 apresenta a configuração inicial do autômato adaptativo equivalente ao autômato de pilha apresentado anteriormente. No estado 1, o autômato recebe os *tokens* e transita para o estado 2 quando processa um sintagma substantivo (SS) ou quando transita em vazio. No estado 2, o autômato transita para si mesmo quando recebe qualquer tipo de verbo: Vi, Vtd, Vlig, Vtpred, Vtdi e Vti. Todas as outras ramificações são criadas por meio de funções adaptativas chamadas em função do tipo de verbo processado.

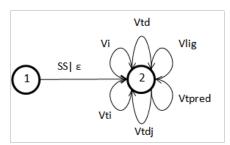


Figura 4.2. Configuração Inicial do Reconhecedor Gramatical.

Por exemplo, se o verbo é de ligação (Vlig), o autômato utiliza as funções adaptativas α (j) e β (o), definidas da seguinte forma:

$$\begin{array}{lll} \alpha \ (j) \colon \{ \ o^* : & \beta \ (o) \colon \{ \ t^*u^*v^*x^* : \\ & \ - \ [\ (j, \ Vlig)] \\ & \ + \ [\ (j, \ Vlig) : \to o, \beta \ (o)] \\ & \ \} & + \ [\ (o, \ Sady) : \to u] \\ & \ + \ [\ (o, \ Sadj) : \to v] \\ & \ + \ [\ (o, \ SS) : \to x] \end{array}$$

A função adaptativa α (j) é chamada pelo autômato antes de processar o *token*, criando o estado 11 e a produção que

leva o autômato do estado 2 ao novo estado criado. Em seguida, o autômato chama a função β (o), criando os estados 12, 13, 14 e 15 e as produções que interligam o estado 11 aos novos estados. A Fig. 4.3 mostra a configuração do autômato após o processamento do verbo de ligação. Neste exemplo, o autômato criou apenas os estados 11, 12, 13, 14 e 15 e as respectivas transições, evitando alocar recursos que seriam necessários para criar o autômato completo, conforme apresentado na Fig. 4.1.

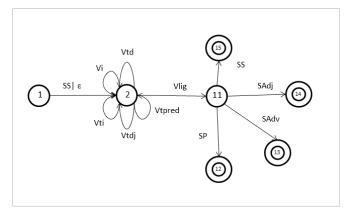


Figura 4.3. Configuração do autômato após o processamento do verbo de ligação.

V. DESAMBIGUAÇÃO SINTÁTICA E SEMÂNTICA

Iwai [20] apresenta um formalismo, denominado Gramáticas Adaptativas, no qual as regras gramaticais são criadas dinamicamente a partir do processamento da cadeia de entrada e de informações relacionadas ao contexto em que se apresentam, permitindo a identificação e resolução de ambiguidades semânticas. A autora demonstra, ainda, a equivalência computacional entre as Gramáticas Adaptativas e os autômatos adaptativos, o que permite que frases da linguagem gerada por uma gramática adaptativa sejam reconhecidas por autômatos adaptativos.

O Linguístico incorpora o formalismo proposto pela autora, através do mapeamento dos sintagmas e padrões oracionais de Luft para uma gramática adaptativa. As ações adaptativas são usadas para modificar as regras de produção em função do contexto definido pelo tipo do verbo encontrado no texto e pela presença de conjunções ou pronomes relativos. O resultado final é o seguinte:

 $G = (G^0, T, R^0)$, onde:

 $G^0 = (VN^0, VT, VC, PL^0, PD^0, S)$

 R^0 = relação que associa regras de produção às ações adaptativas.

 $VN^0 = \{S, con, SS, SV, SS, Sadj, SP, num, PrA, Sc, Sp,$ PrPes, Neg, Aux, PreV,V, Vlig, Vtda,Vtdna, Vtdi, Vtpred ,Vti ,Vi, Sadv, Prep, Ind, ArtDef, ArtInd, Dem, Pos, PrRel, Adv, Adj, pont }

 $VT = \{\{\{1 \le i \le 0\}, \{,\}, ;, ., :, !, ?, 1, 2, 3, 4, 5, 6, 7, 8, 9, 0\}\}$

 $VC = \emptyset$

 $PD^0 = \emptyset$

 $PL^0 = \{ S \rightarrow [con] [PrRel] [SS] SV pont [S] \}$

 $SS \rightarrow [Sadj] SS [Sadj | SP] [S] | S$

 $SS \rightarrow (([num \mid PrA] Sc) \mid Sp \mid PrPes)$

 $SV \rightarrow [Neg] [Aux \mid PreV] V [SS \mid Sadj \mid Sadv \mid SP] \mid ((SS \mid Sadj \mid Sadv \mid SP)) \mid ((SS \mid Sadv \mid SAVv \mid SP)) \mid ((SS \mid Sadv \mid SAVv \mid SAVv$ SP) | (SS Sadj) | (SP SP) | (SS SS) | (SS Sadv) | (SS SP SP))

 $V \rightarrow (Vlig \mid Vtda \mid Vtdna \mid Vtdi \mid Vtpred \mid Vti \mid Vi)$

 $SP \rightarrow Prep (SS \mid Sadj)$

 $Sadj \rightarrow Sadj [SP] [S] | S$

 $Sadj \rightarrow [Adv] Adj [S]$

 $Sadv \rightarrow Sadv [SP] [S] | S$

 $Sadv \rightarrow [Adv] Adv [S]$

 $PrA \rightarrow [PrA] (Ind \mid ArtDef \mid ArtInd \mid Dem \mid Pos)$

con → (aditiva | adversativa | alternativa | conclusiva| explicativa | integrante | causal | comparativa | concessiva | condicional | conformativa | consecutiva | final | proporcional | temporal)

multiplicativo| fracionario)

 $pont \rightarrow (, |; |. |: |! |?)$

Sendo:

S – Sentença

SS – Sintagma substantivo

SV - Sintagma verbal

SP - Sintagma preposicional

SN – Sintagma nominal

Sadv – Sintagma adverbial

Sadj – Sintagma adjetivo

Adv – adverbio

Adj – adjetivo

ArtDef – artigo definido

ArtInd - artigo indefinido

Aux – Partícula auxiliar (apassivadora ou pré-verbal)

Conec – Conector (conjunção ou pronome relativo)

Dem – pronome demonstrativo indefinido

Pos – pronome possessivo

Ind – pronome indefinido

Neg – partícula (negação)

PrA – pronome adjetivo

PrPes – pronome pessoal

Prep – preposição

PreV - pré-verbais

Quant – numeral

Sc – substantivo comum

Sp – substantivo próprio

V– verbo

Vlig – verbo de ligação

Vi – verbo intransitivo

Vtd - verbo transitivo direto

Vti – verbo transitivo indireto

Vtdi – verbo transitivo direto e indireto

A análise semântica das sentenças é feita de acordo com os termos processados, que podem gerar novas regras de produção e exclusão de regras existentes. Por exemplo, quando o Linguístico encontra um verbo transitivo indireto, ele dispara uma ação adaptativa A (Vti) que elimina regras iniciais definidas na gramática G⁰ e cria outras, relacionadas aos complementos oracionais exigidos por este tipo de verbo.

A representação destas alterações é feita da seguinte forma:

$$\begin{split} \text{A} \quad (\text{ } t = Vti) = \text{L} \quad (t) = \{ \text{ } + [\text{R: }_{\text{G}}{}^{0}\text{P26-66, 69,77-78, 87, 93, 101-102, 113-114:} \\ & \text{G}^{1}\text{V} \leftarrow \chi \backslash \text{V}] \\ & + [\text{G1}\chi\text{V} \rightarrow \text{Vti}] \\ & - [\text{R: } \text{G}^{0}\text{P67-68, 70-76, 79-86, 88-100, 103-112, 115-120: } \text{G}^{1}\varnothing] \end{split}$$

A sintaxe da instrução R é a seguinte:

 $R: <_G^n > < P_{I-F} > : <_G^{n+1} > P'$

onde:

R: replace

_Gⁿ: gramática anterior à atualização

 $P_{I\text{-}F}$: regras de produção para atualização, sendo I= Inicio de intervalo e F= fim de intervalo. Acrescenta-se mais de um grupo de regras usando a vírgula para separação.

Gⁿ⁺¹: gramática posterior à atualização

P': nova regra de produção

 \varnothing : símbolo que indica que as regras assinaladas em $P_{\text{I-F}}$ serão removidas

 $V \leftarrow \chi \backslash V$: indica a aplicação de contexto ao não-terminal V

χV→Vti: indica a derivação de um não terminal V, dependente de contexto, em um não terminal Vti

O formalismo apresentado por Iwai é bastante prático e fornece suporte para incluir características adicionais que o parser pode utilizar para resolver problemas em que o contexto seja importante para definição das regras de produção. Por exemplo, o parser pode verificar se a frase analisada está correta ou resolver possíveis ambiguidades caso tenha à disposição informações sobre o tipo do verbo, tempo, modo e pessoa do verbo analisado. No exemplo abaixo, a ação adaptativa A utiliza como parâmetros de entrada a identificação do tipo do verbo, modo e pessoa para definição das funções adaptativas.

$$\begin{split} T &= \{ \text{ A (verbo=Vlig, tempo = Presente, modo = Singular, pessoa = Primeira) = F (verbo, tempo, modo, pessoa) = } \\ \{ + [\text{R: }_G{}^n \text{ P}_{\text{I-F}}\text{: }_G{}^{n+1}\text{V} \leftarrow \chi \text{V}] \\ + [\text{G}^{n+1}\chi\text{V} \rightarrow \text{Vlig, Presente, Singular, Primeira}] \\ + [\text{R: }_G{}^n \text{ P}_{\text{I-F}}\text{: }_G{}^{n+1}\varnothing] \ \} \end{split}$$

Analogamente, é possível incluir regras para análise da concordância nominal ou qualquer outro tipo de conteúdo semântico, sem necessidade de usar qualquer outro elemento fora do formalismo. As gramáticas adaptativas também podem ser usadas para representar o uso de informações probabilísticas na seleção das regras de produção. Tal técnica é usada quando existe mais de uma regra de produção aplicável e não existem informações sintáticas e semânticas suficientes para desambiguá-las. Neste caso, é possível usar a probabilidade de ocorrência das regras como fator de escolha e o formalismo de Iwai pode ser usado representar este tipo de informação contextual. No exemplo abaixo, a ação adaptativa A utiliza como parâmetro de entrada a probabilidade de ocorrência das regras de produção avaliadas e uma indicação para usar a regra de máxima probabilidade.

$$SVprob \rightarrow V SP 10\%$$

 $SVprob \rightarrow V (SS SP) 90\%$

$$\begin{split} T &= \{\; A \; (\; t\text{=SV}, \, u = prob, \, v\text{=max}) = F \; (SV, \, prob, \, max) = \\ \{\; + [R:\;_{G}^{\; n} \; P_{1\text{-}F}:\;_{G}^{\; n+1} \; SV \leftarrow \chi SV prob] \\ + [_{G}^{\; n+1} \; \chi SV_{prob} \rightarrow \chi SV_{max}] \\ ? \; [_{G}^{\; n+1} \; \chi SV_{max}] \\ + [_{G}^{\; n+1} \; \chi SV_{max} \rightarrow V \; (SS \; SP)] \\ + [R:\;_{G}^{\; n} \; P_{1\text{-}F}:\;_{G}^{\; n+1} \; \varnothing] \; \} \end{split}$$

VI. EXPERIMENTOS E RESULTADOS COMPARATIVOS

Os módulos apresentados na seção IV foram desenvolvidos em linguagem Python, com o apoio do NLTK [21], uma biblioteca de processamento de linguagem natural. O Sentenciador e o *Token*izador foram implementados exatamente como haviam sido descritos previamente. Já o Identificador Morfológico e o Analisador Sintático foram implementados a partir do corpus CINTIL-Treebank [22] e das regras de produção da Teoria X Barra [23] na qual ele se apoia. O CINTIL-Treebank, corpus da Língua Portuguesa desenvolvido pela Universidade de Lisboa, foi escolhido nesta etapa da pesquisa por usar o mesmo sistema de notação da Universidade da Pennsylvania - Treebank, que é considerado padrão no processamento de Linguagem Natural, e por se apoiar em uma gramática estrutural e não descritiva, assim como a de Luft, fornecendo padrões oracionais equivalentes, sem perda efetiva do poder de análise. A mudança permitiu a realização de testes que não seriam possíveis com a gramática de Luft, pois ela não está disponível em corpus pré-anotados e, ainda, a comparação dos resultados de trabalhos publicados em diferentes línguas.

O algoritmo Viterbi descrito na seção IV, Fig. 2.10, foi implementado através de um mecanismo de *fall-back*, analisando trigramas, bigramas e unigramas, e assumindo a classificação mais comum do corpus, que é substantivo, caso não tenha encontrado a classificação do *token* analisado. O treino utilizou 90% do corpus CINTIL-Treebank e os testes foram feitos em 10%, escolhidos aleatoriamente. A Tabela 4 apresenta uma comparação dos resultados do Linguístico e dos analisadores morfológicos do estado da arte [24,25].

TABELA 4. COMPARAÇÃO DE RESULTADOS

Sistema	Linguístico	TBL	TnT	MXPost	QTag
Precisão	91,41%	97,09%	96,87%	97,08%	89,97%

Observa-se que os resultados obtidos estão abaixo daqueles apresentados por analisadores morfológicos do estado da arte, que chegam a 97,09% de precisão. No entanto, é possível obter melhorias, visto que foi usado um corpus de treinamento de tamanho reduzido e também por não terem sido utilizados recursos semânticos e contextuais de desambiguação, que serão incorporados nos próximos experimentos.

O Reconhecedor Sintático utilizou as regras de produção capturadas do CINTIL-Treebank, enriquecidas com as probabilidades de ocorrências obtidas a partir de um subconjunto do mesmo corpus. As regras de produção foram extraídas diretamente do CINTIL-Treebank, usando bibliotecas do NLTK e realizando algumas adaptações, visando evitar a repetição de regras e incorporar o *encoding* usado na Língua Portuguesa (ISO-Latin-1). A Fig.5 apresenta um fragmento das regras obtidas.

AP -> A CONJP [0.00766]
ADV> PNT ADV_ PNT [0.05036]
S -> V_ AP [0.00019]
$AD \sim ADW = A \cdot 0.007661$

VP -> V ADV [0.01065] A_ -> ADV A_ [0.01639] AP -> AP CP [0.00128] NP -> NP QNT [0.00003]

Figura 5 – Fragmento das Regras de Produção Extraídas do CINTIL – Treebank.

As regras de produção foram estruturadas de forma a apresentar as derivações dos não terminais em outros não terminais, seguidos da probabilidade calculada para a regra, em parênteses. Usando, como exemplo, o fragmento da Fig.5, o não terminal AP deriva para um não terminal A e outro não terminal CONJP com probabilidade de 0.00766 ou 0.766 %. As derivações de não terminais para terminais foram excluídas da gramática, sendo substituídas pelas classificações geradas pelo Identificador Morfológico.

Os resultados foram medidos utilizando o método Parseval [26], considerado padrão para apurar o grau de correção das árvores de derivação geradas por parsers. Os testes foram feitos com 1070 frases do corpus CINTIL-Treebank, escolhidas aleatoriamente, representando aproximadamente 10% do total. A Tabela 5 apresenta os resultados obtidos pelo Reconhecedor Sintático, comparando-o a outros parsers. O índice FParseval (média harmônica de precisão e cobertura) foi de 91,49%, melhor, portanto, que os resultados apresentados pelos parsers Bikel, Stanford e Berkeley aplicados ao mesmo corpus - 84,97%, 88,07% e 89,33%, respectivamente, chegando a níveis semelhantes aos obtidos com a língua inglesa [22]. Os resultados, no entanto, precisam ser relativizados pelo fato de tanto o treino quanto o teste terem sido realizados em um mesmo corpus, o que pode ter gerado algum tipo de viés, mesmo dividindo o corpus em subconjuntos diferentes para cada atividade.

TABELA 5. COMPARAÇÃO DE RESULTADOS

Sistema	Linguístico	Bikel	Stanford	Berkeley
FParseval	91,49 %	84,97%	88,07%	89,33%,

VII. CONSIDERAÇÕES FINAIS

Este artigo apresentou uma revisão dos conceitos de Tecnologia Adaptativa e de Processamento da Linguagem Natural. Em seguida, foi apresentado o Linguístico, uma proposta de reconhecedor gramatical que utiliza autômatos adaptativos como tecnologia subjacente. Procurou-se também apresentar cenários em que a gramática utilizada pelo Linguístico é estendida, incorporando critérios semânticos para auxiliar na resolução de ambiguidades. Ao final, foram apresentados os experimentos realizados com o Linguístico, comparando os resultados com os do estado da arte.

O Identificador Morfológico apresentou resultados um pouco abaixo dos analisadores morfológicos de referência. Dado que foram os primeiros experimentos e que não houve uso de informações semânticas e contextuais para desambiguação, espera-se melhora com o desenvolvimento da pesquisa. Já o Reconhecedor Sintático apresentou resultados melhores do que os *parsers* do estado da arte aplicados a textos em português e chegou a níveis semelhantes aos apresentados em análises de textos em inglês.

Os resultados observados nos encorajam a avançar na investigação, testando o Linguístico com outros corpus e implementando melhorias já identificadas nesta etapa do trabalho. Espera-se também ampliar o escopo da pesquisa, com a inclusão de aspectos de dependência do contexto, equivalências entre gramáticas, dialetos e regionalismos, estendendo o uso das técnicas e métodos adaptativos no processamento de linguagem natural.

REFERÊNCIAS

- [1] http://www.pcs.usp.br/~lta/
- [2] Taniwaki, C. Formalismos adaptativos na análise sintática de Linguagem Natural. Dissertação de Mestrado, EPUSP, São Paulo, 2001
- [3] Menezes, C. E. Um método para a construção de analisadores morfológicos, aplicado à língua portuguesa, baseado em autômatos adaptativos. Dissertação de Mestrado, Escola Politécnica da Universidade de São Paulo, 2000.
- [4] Padovani, D. Uma proposta de autômato adaptativo para reconhecimento de anáforas pronominais segundo algoritmo de Mitkov. Workshop de Tecnologias Adaptativas – WTA 2009, 2009.
- [5] Moraes, M. Alguns aspectos de tratamento sintático de dependência de contexto em linguagem natural empregando tecnologia adaptativa, Tese de Doutorado, Escola Politécnica da Universidade de São Paulo, 2006.
- [6] Rich, E.; Knight, K. Inteligência Artificial, 2. Ed. São Paulo: Makron Books, 1993.
- [7] Vieira, R.; Lima, V. Linguística computacional: princípios e aplicações. IX Escola de Informática da SBC-Sul, 2001.
- [8] Fuchs, C.; Le Goffic, P. Les Linguistiques Contemporaines.
- [9] M. G. V.Nunes et al. Introdução ao Processamento das Línguas Naturais. Notas didáticas do ICMC Nº 38, São Carlos, 88p, 1999. Paris, Hachette, 1992. 158p.
- [10] Sardinha, T. B. A Língua Portuguesa no Computador. 295p. Mercado de Letras, 2005.
- [11] Rocha, R.L.A. Tecnologia Adaptativa Aplicada ao Processamento Computacional de Língua Natural. Workshop de Tecnologias Adaptativas – WTA 2007, 2007.
- [12] Luft, C. Moderna Gramática Brasileira. 2ª. Edição Revista e Atualizada. 265p. Editora Globo, 2002.
- [13] http://www.linguateca.pt/
- [14] http://www.nilc.icmc.usp.br/tep2/
- [15] http://www.portalsaofrancisco.com.br/alfa/materias/index-linguaportuguesa.php/
- [16] Basilio, M. Formação e Classes de Palavras no Português do Brasil. Ed.Contexto, 2004.
- [17] Basilio, M. Teoria Lexical. Ed.Atica, 1987.
- [18] http://www.cs.columbia.edu/~mcollins/hmms-spring2013.pdf/
- [19] Forney, G.D. J. IEEE. Proceedings of the IEEE, Volume 61, Issue 3, 1973.
- [20] Iwai, M.: Um Formalismo Gramatical para Linguagens Dependentes de Contexto. São Paulo, 2000. Escola Politécnica da Universidade de São Paulo. Tese de Doutorado (2000).
- [21] Bird, Steven, Edward Loper and Ewan Klein (2009), Natural Language Processing with Python. O'Reilly Media Inc.
- Branco, António, Francisco Costa, João Silva, Sara Silveira, Sérgio Castro, Mariana Avelãs, Clara Pinto and João Graça, 2010, "Developing a Deep Linguistic Databank Supporting a Collection of Treebanks: the CINTIL DeepGramBank ", In Proceedings, LREC2010 The 7th international conference on Language Resources and Evaluation, La Valleta, Malta, May 19-21, 2010.
- [23] Mioto C., Silva M. C. F., Lopes, R. Novo Manual de Sintaxe. Ed. Contexto. 2013.
- Branco, A., Silva, J.: Evaluating solutions for the rapid development of state-of-the-art POS taggers for Portuguese. In: Proceedings of the 4th Language Resources and Evaluation Conference (LREC). (2004) 507–510.
- [25] Silva, J.: Shallow processing of Portuguese: From sentence chunking to nominal lemmatization. Master's thesis, University of Lisbon (2007) Published as Technical Report DI-FCUL-TR-07-16.

[26] Black, E., Abney, S., Flickinger, D., Gdaniec, C., Grishman, R., Harrison, P., Hindle, D., Marcus, M., Santorini, B.: A procedure for quantitatively comparing the syntactic coverage of English grammars. In: Proceedings of the Workshop on the Evaluation of Parsing Systems. (1991) 306–311.

Ana Teresa Contier formou-se em Letras-Português pela Universidade de São Paulo (2001) e em publicidade pela PUC-SP (2002). Em 2007 obteve o título de mestre pela Poli-USP com a dissertação: "Um modelo de extração de propriedades de textos usando pensamento narrativo e paradigmático".

Djalma Padovani formou-se em administração de empresas pela Faculdade de Economia e Administração da Universidade de São Paulo, em 1987 e obteve o mestrado em engenharia de software pelo Instituto de Pesquisas Tecnológicas de São Paulo - IPT, em 2008. Trabalhou em diversas empresas nas áreas de desenvolvimento de software e tecnologia de informação e atualmente atua no Laboratório de Dados da Serasa Experían.

João José Neto graduado em Engenharia de Eletricidade (1971), mestrado em Engenharia Elétrica (1975) e doutorado em Engenharia Elétrica (1980), e livre-docência (1993) pela Escola Politécnica da Universidade de São Paulo. Atualmente é professor associado da Escola Politécnica da Universidade de São Paulo, e coordena o LTA - Laboratório de Linguagens e Tecnologia Adaptativa do PCS - Departamento de Engenharia de Computação e Sistemas Digitais da EPUSP. Tem experiência na área de Ciência da Computação, com ênfase nos Fundamentos da Engenharia da Computação, atuando principalmente nos seguintes temas: dispositivos adaptativos, tecnologia adaptativa, autômatos adaptativos, e em suas aplicações à Engenharia de Computação, particularmente em sistemas de tomada de decisão adaptativa, naílise e processamento de linguagens naturais, construção de compiladores, robótica, ensino assistido por computador, modelagem de sistemas inteligentes, processos de aprendizagem automática e inferências baseadas em tecnologia adaptativa.