

**Universidade de São Paulo**  
**Escola de Artes, Ciências e Humanidades - EACH**

Projeto CAPES 6415/10-5

**Extensão de Algoritmos de Geração de Expressões  
de Referência para Domínios Físicos**

Ivandr  Paraboni (USP / EACH)

## Resumo

A geração de expressões de referência (GER) é um dos componentes fundamentais de aplicações de geração de língua natural (GLN) a partir de dados de entrada não linguísticos. Do ponto de vista computacional, a tarefa de GER tem sido tratada considerando-se principalmente domínios (físicos ou discursivos) simplificados e/ou bidimensionais. O problema de referência em domínios físicos mais realistas (e.g., com grande complexidade estrutural e/ou tridimensionalidade), entretanto, permanece praticamente inexplorado na pesquisa da área. Neste trabalho propomos uma extensão de algoritmos de GER (desenvolvidos em nosso trabalho prévio) para domínios físicos tais como ambientes virtuais tridimensionais, e sua validação empírica em colaboração com um grupo de pesquisa estrangeiro de grande tradição na área.

**Palavras-chave:** processamento de língua natural, geração de língua natural, expressões de referência.

## Sumário

1	Introdução.....	3
2	O problema computacional de GER.....	5
3	Proposta de Pesquisa.....	8
3.1.	Motivação.....	8
3.2.	Objetivos.....	9
3.3.	Metodologia.....	9
3.4.	Atividades Previstas.....	10
3.5.	Colaboração externa.....	11
3.6.	Resultados Esperados.....	13
4	Considerações Finais.....	13

## 1 Introdução

Considere a cena a seguir, na qual desejamos apontar uma pessoa em particular, ou seja, *fazer referência* a um determinado indivíduo:



**Fig. 1.** Um ambiente virtual tridimensional com elementos referenciáveis.

Uma referência a uma pessoa deste exemplo poderia ser feita de várias formas, e em geral sem muito esforço: participantes (humanos) de um discurso facilmente produziram expressões de referência como "o rapaz de braços cruzados, à esquerda", "o homem de blusão escuro" ou "a moça mais à frente".

Considere agora esta mesma tarefa sendo desempenhada por um agente de software em um ambiente envolvendo comunicação humano-computador, ou seja, em uma aplicação típica de Geração de Língua Natural<sup>1</sup> (GLN). Dado um objeto-alvo ao qual deseja-se fazer referência (por exemplo, uma pessoa) e uma base de dados contendo todos os fatos relevantes sobre os elementos do domínio (e.g., propriedades semânticas como cores, tamanhos, distâncias, posições etc.), o objetivo de um algoritmo de *geração de expressões de referência* (GER) é produzir uma lista de propriedades do objeto-alvo tal que esta seja capaz de distingui-lo de todos os demais, constituindo uma *descrição ótima* para aquela situação.

---

<sup>1</sup> Sistemas deste tipo - que produzem descrições textuais a partir de uma entrada de dados não linguística - são empregados quando texto predefinido não é suficiente, ou seja, quando é necessária uma variação e sofisticação linguística comparáveis ao desempenho humano, e são comuns em aplicações de cunho educativo, comercial, médico, de entretenimento e outras.

Decidir o que exatamente constitui uma descrição ótima, ou quais propriedades do objeto-alvo ela deve conter, é o problema central para um algoritmo de GER. Entre os muitos objetivos conflitantes a considerar nesta tomada de decisões, temos por exemplo que a descrição resultante não deve ser *ambígua* (e.g., "o rapaz"), nem excessivamente *breve* (e.g., "a de verde", em referência à moça mais ao fundo), ou excessivamente *longa* (e.g., "o rapaz de braços cruzados, à esquerda, e que está ao lado do homem de blusão escuro").

Além disso, a seleção de conteúdo deve considerar também a *saliência relativa* das propriedades disponíveis. Por exemplo, a descrição "o homem com os braços à cabeça" pode ser preferível à "o homem de calça marrom", que mesmo sendo uma descrição não ambígua pode exigir maior esforço de interpretação. Ou, de forma mais ampla, podemos afirmar que algoritmos de GER preocupam-se em produzir descrições *psicologicamente plausíveis*, ou seja, tão próximas quanto possível das que seriam produzidas por um agente humano sob as mesmas circunstâncias.

Estes exemplos ilustram apenas alguns dos muitos desafios que caracterizam o problema computacional de GER. Apesar do escopo à primeira vista limitado, a geração de expressões de referência é na verdade uma questão central na área de GLN, e não por acaso uma de suas linhas de pesquisa mais ativas (e.g., Dale & Reiter, 1995; Krahmer & Theune, 2002; Krahmer et. al., 2003; Horacek, 2005; van Deemter et. al., 2006; Gatt et. al., 2007, 2008, 2009; Turner et. al., 2008, 2009; Dale & Viethen, 2009; Mitchell et. al., 2010).

O problema computacional de GER é também a principal linha de pesquisa do proponente deste projeto, iniciada no período de doutoramento (Paraboni, 2003; Paraboni et. al., 2006, 2007) e complementada em diversas oportunidades com o apoio de alunos de Iniciação Científica desta instituição (Lucena et. al., 2010; Lucena & Paraboni, 2008, 2008a, 2009, 2009a, 2010). O desenvolvimento de algoritmos de GER foi também abordado em um projeto de pesquisa individual recente<sup>2</sup>, e é o foco da presente proposta.

---

<sup>2</sup> CNPq edital Universal 15/2007, nro. 484015/2007-9, concluído em dezembro/2009.

## 2 O problema computacional de GER

Uma definição clássica do problema de seleção de conteúdo para GER é encontrada na especificação de um dos mais bem conhecidos e influentes algoritmos da área, o *algoritmo Incremental* (Dale & Reiter, 1995). Este algoritmo recebe como entrada um objeto-alvo ou referente  $r$  que se deseja descrever, um contexto formado pelos objetos que são passíveis de referência em um dado ponto do discurso, e suas propriedades semânticas referenciáveis (i.e., aquelas que podem figurar em uma descrição definida). O objetivo do algoritmo é computar uma lista  $L$  de propriedades semânticas de  $r$ , tal que  $L$  permita a identificação de  $r$  sem ambiguidade, i.e., tal que  $L$  possa ser realizada na forma de uma expressão linguística que descreva apenas o objeto-alvo, e nenhum outro objeto do contexto.

As propriedades do objeto-alvo  $r$  são incluídas em  $L$  de forma sequencial segundo uma ordem de preferência  $P$  determinada pelo domínio, até que uma descrição única (i.e., livre de ambiguidade) seja obtida. As propriedades consideradas para inclusão são limitadas àquelas capazes de distinguir  $r$  dos demais objetos no contexto, ou seja, o algoritmo considera apenas a inclusão de propriedades *restritivas* ou sem conteúdo redundante (do ponto de vista lógico).

Por exemplo, suponha um contexto com quatro pessoas como na Fig.1 anterior, e propriedades referenciáveis organizadas em uma base de dados como segue:

- $e_1$ : gênero=m, braços=cruzados, cabelo=claro, direção=costas.
- $e_2$ : gênero=f, braços=cruzados, cabelo=escuro, direção=esquerda.
- $e_3$ : gênero=m, braços=à cabeça, cabelo=escuro, direção=frente.
- $e_4$ : gênero=f, braços=cruzados, cabelo=escuro, direção=direita.

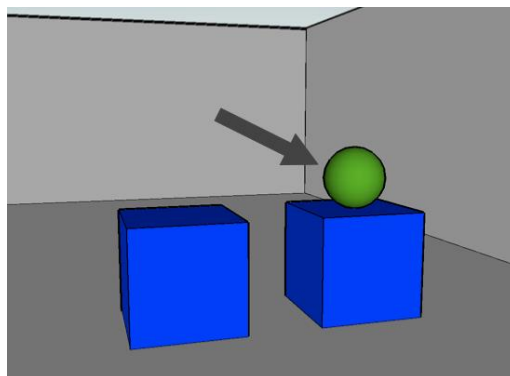
Suponha que a ordem de preferência para seleção de conteúdo seja  $P = \{\textit{gênero}, \textit{braços}, \textit{cabelo}, \textit{direção}\}$ , e que o objeto-alvo em questão seja  $r = e_2$ . Neste caso, o algoritmo Incremental consideraria primeiro o atributo *gênero*, o qual seria incluído na descrição resultante  $L$  porque seu valor "f" elimina dois elementos do contexto ( $e_1$  e  $e_3$ , cujo *gênero* possui o valor "m"). A seguir, o atributo *braços* é considerado para inclusão, mas este não será selecionado porque seu valor "cruzados" não descarta o único elemento restante no contexto ( $e_4$ ), que também possui este mesmo valor. O atributo *cabelo* é desconsiderado pelo mesmo motivo, e finalmente o atributo *direção* é incluído porque seu valor "esquerda" descarta  $e_4$  (cuja *direção* possui o valor "direita").

Assim, a expressão resultante  $L = \{\text{gênero}=f, \text{direção}=esquerda\}$  descreve de forma única o objeto-alvo, e poderia ser realizada, por exemplo, como "a moça voltada para a esquerda". De forma análoga,  $e_1$  poderia ser descrito como "o rapaz de braços cruzados",  $e_3$  poderia ser descrito como "o rapaz com os braços à cabeça" e  $e_4$  como "a moça voltada para a direita".

A preferência por propriedades discriminatórias no algoritmo Incremental é motivada pela máxima de *brevidade* proposta por H. P. Grice (Grice, 1975), com o objetivo de minimizar possíveis implicações lógicas que possam frustrar o significado comunicado. Por exemplo, em "o rapaz de braços à cabeça, próximo de uma cadeira" a referência à posição física é redundante, e em certos contextos pode ser potencialmente prejudicial à compreensão do sentido comunicado, já que pode sugerir a existência de um segundo rapaz, também com as mãos à cabeça, e que presumivelmente não estaria próximo de uma cadeira.

No algoritmo Incremental e abordagens relacionadas, os objetos do domínio são descritos em termos de suas propriedades *atômicas* (e.g., tipo, cor, tamanho etc.). Nos poucos casos em que propriedades *relacionais* são consideradas (e.g., Dale & Haddock, 1991; Krahmer & Theune, 2002; Dale & Viethen, 2009), os exemplos discutidos são simplificados de forma que a dificuldade de identificação (e.g., a busca pelo objeto-alvo da expressão de referência) não é uma questão relevante.

Considere o exemplo a seguir, ilustrando o domínio de GER abordado em Dale & Viethen (2009), e que constitui um dos poucos estudos do gênero a tratar algum tipo de tridimensionalidade, neste caso envolvendo referência espacial entre objetos geométricos esquemáticos.



**Fig. 2.** Domínio tridimensional simplificado considerado em Dale & Viethen (2009).

O ponto que gostaríamos de chamar à atenção na presente proposta é que, apesar de representar um avanço considerável em relação ao tratamento de domínios puramente bidimensionais, estudos como o de Dale & Viethen (2009) ainda partem do pressuposto de que *todos os fatos do domínio são acessíveis aos participantes do discurso*, e isso, entretanto, é *claramente insuficiente* no caso de um domínio físico mais complexo, como o mundo virtual da Fig.1 anterior.

Para um algoritmo de GER, a diferença entre os domínios da Fig. 1 e 2 (ou seja, a diferença entre assumir que os fatos do domínio são parcialmente ou totalmente acessíveis) se reflete essencialmente na necessidade de tratamento da *informação redundante* para localização de referentes. Por exemplo, no contexto da Fig.1 anterior, dada o objetivo de produzir uma descrição única para o homem mais à esquerda (que na ilustração colorida é a única pessoa de cabelos claros na cena), o algoritmo Incremental e abordagens derivadas produziria algo como “o homem loiro”, que é uma forma breve e não ambígua de descrever aquele objeto-alvo.

No entanto, se o grupo de pessoas em questão fosse mais numeroso, uma expressão deste tipo – embora livre de ambiguidade - poderia ser de difícil resolução. Neste caso, seria mais adequado incluir algum tipo de informação redundante, como em “o homem loiro, *de blusa azul*”, mesmo não sendo esta informação estritamente necessária para fins de identificação do objeto-alvo. De forma similar, o algoritmo de GER precisaria incluir informações não previstas pelas abordagens tradicionais se, por exemplo, o objeto-alvo estivesse oculto por outros (e.g., atrás de um vaso de plantas) ou em uma posição mais afastada.

Assim, para gerar expressões de referência em um domínio bidimensional, ou mesmo um domínio tridimensional simplificado como os cubos e esferas da Fig.2, provavelmente bastaria empregar uma versão relacional do algoritmo de Dale & Reiter, como apresentada em Krahmer & Theune (2002) e outros. Já para situações em que os fatos sobre o domínio não são (do ponto de vista humano) instantaneamente reconhecíveis, abordagens de GER tradicionais não produziriam descrições de utilidade prática, fazendo-se necessário fornecer informações adicionais (e que do ponto de vista lógico seriam redundantes) para viabilizar a identificação do objeto-alvo.

Algoritmos com estas características (i.e., que não pressupõe conhecimento total dos fatos do domínio) foram desenvolvidos pelo proponente desta pesquisa em duas ocasiões distintas: em Paraboni et. al. (2007) foi proposta uma série de

algoritmos de GER que facilitam a busca por referentes em contextos físicos do tipo *hierárquico* formando partições (e.g., países divididos em estados e cidades, um campus dividido em prédios e salas etc.), contemplando três estratégias distintas para atingir o equilíbrio entre facilidade de busca (favorecendo descrições mais extensas) e facilidade de interpretação (favorecendo descrições mais breves).

Posteriormente, este estudo foi estendido para o domínio discursivo (Lucena et. al., 2010), investigando-se o uso de informação redundante em referências anafóricas. Neste e em trabalhos correlatos (Lucena & Paraboni, 2008, 2008a, 2009, 2009a, 2010) foram propostos novos algoritmos de GER explorando uma estratégia geral que combina a seleção de atributos altamente discriminatórios (favorecendo descrições breves) e a seleção de atributos altamente frequentes (favorecendo descrições mais plausíveis).

Nenhum dos algoritmos propostos em nossos trabalhos prévios, entretanto, contempla situações de referência como no domínio físico exemplificado na Fig.1. Uma proposta de pesquisa desta natureza é apresentada nas seções a seguir.

### **3 Proposta de Pesquisa**

#### **3.1. Motivação**

Em domínios físicos do tipo considerado nesta proposta, um objeto-alvo poderia ter não apenas um grande número de objetos concorrentes no mesmo foco de atenção, mas estes objetos (bem como suas propriedades) teriam diferentes graus de saliência relativa. Além disso, certos objetos poderiam estar parcial ou totalmente ocultos por outros (e.g., dentro de uma caixa), ou pela própria estrutura do domínio (e.g, atrás de uma parede), dentre muitas outras possibilidades.

Modelar todas estas características de forma que sejam levadas em conta por um algoritmo de GER não é uma tarefa trivial, e a geração de expressões de referência de forma irrestrita em ambientes como o da Fig.1. anterior é um cenário ainda bastante distante do estado da arte. Como uma aproximação inicial do problema, entretanto, podemos tirar proveito de trabalho prévio enfocando tanto a referência em domínios físicos hierárquicos (Paraboni et. al., 2007) quanto discursivos (Lucena et. al., 2010). Nestes dois casos, embora os fenômenos tratados sejam distintos do tema da presente discussão, foram propostos algoritmos que levam em conta o uso



de informação redundante na identificação do objeto-alvo referenciado. E este princípio, conforme sugerido na seção anterior, é aplicável (ou pelo menos adaptável) aos tipos de domínios discutidos nesta proposta.

A aplicação de algoritmos de GER que fazem uso de informação redundante (e.g., Paraboni et. al., 2007; Lucena et. al., 2010) em um domínio físico mais realista é, até onde temos conhecimento, uma iniciativa de pesquisa inédita em GLN/GER, e que vem preencher uma clara lacuna no conhecimento da área. Além disso, uma investigação desta natureza nos parece perfeitamente viável uma vez que há um volume considerável de trabalho teórico já realizado, e uma série de algoritmos básicos já disponibilizados.

### **3.2. Objetivos**

O objetivo da presente proposta de pesquisa é a extensão de algoritmos de seleção de conteúdo semântico desenvolvidos em trabalho prévio (e.g., Paraboni et. al., 2007; Lucena et. al., 2010) para *domínios físicos*, e sua validação através de experimentos práticos de GER.

### **3.3. Metodologia**

O projeto consiste na adaptação de algoritmos já existentes para produzir referências em domínios físicos (como o ambiente virtual 3D utilizado como exemplo na introdução desta proposta) e na condução de experimentos para sua validação em colaboração com especialistas vinculados a uma instituição de pesquisa no exterior.

Considerando-se que uma parte do trabalho (em especial, os algoritmos a serem tomados como base) já se encontra desenvolvida, o projeto proposto é de curta duração (6 meses), contemplando apenas as atividades a serem realizadas idealmente em colaboração. De forma mais específica, o projeto privilegia a etapa de validação empírica de algoritmos de GER, área esta na qual a interação com o grupo no exterior faz-se mais desejável em razão da sua experiência em pesquisas deste tipo (e.g., van Deemter et. al., 2006; Gatt et. al., 2007, 2008, 2009; Turner et. al., 2008, 2009; Mitchell et. al., 2010) conforme será discutido a seguir.

### **3.4. Atividades Previstas**

#### *1-Revisão Bibliográfica*

Estudo de algoritmos tradicionais como Dale & Reiter (1995) e Krahmer et. al. (2003) para seleção de conteúdo de expressões de referência, bem como abordagens recentes enfocando o fenômeno de referência em domínios físicos (e.g., Dale & Viethen, 2009, Mitchell et. al., 2010).

#### *2-Refinamento dos algoritmos propostos*

Os algoritmos propostos em Paraboni et. al. (2007) e Lucena et. al. (2010) serão apresentados ao grupo de pesquisa da instituição de destino para discussão. Serão levantadas as questões de pesquisa mais pertinentes considerando-se o propósito de referência em domínios físicos, e sua especificação será adaptada, em especial, para levar em conta a questão da saliência relativa dos objetos do domínio, nos moldes de Krahmer & Theune (2002). O resultado desta atividade é o refinamento de uma proposta de algoritmo de GER para domínios físicos a ser validada empiricamente.

#### *3-Implementação*

Desenvolvimento da proposta de algoritmo refinada, a partir de implementações já existentes (Paraboni, 2003; Lucena & Paraboni, 2008, 2008a).

#### *4-Definição de experimentos*

Detalhamento de experimentos de GER para validação do algoritmo selecionado e sua comparação com algoritmos de *baseline* que sejam pertinentes. Em especial, consideramos para fins de comparação o algoritmo Incremental (Dale & Reiter, 1995), os algoritmos para identificação de referentes em domínios estruturalmente complexos de Paraboni et. al., (2007) e os algoritmos que combinam estratégia de seleção gulosa e propriedades frequentes (Lucena & Paraboni, 2008).

Os experimentos consistem de testes dos algoritmos selecionados em um ambiente físico possivelmente simulado através da plataforma GIVE (Byron et. al., 2007) amplamente utilizada na instituição de destino. Neste tipo de ambiente, as reações de usuários (tempo de interpretação e resolução de expressões, número de

erros cometidos, etc.) diante de expressões produzidas por algoritmos variados serão monitorada para fins de validação nos moldes de Paraboni et. al. (2007).

### *5-Execução*

Execução dos experimentos previstos e possível revisão e/ou complementação destes. Nesta etapa em especial esperamos poder tirar máximo proveito da experiência do grupo de pesquisa estrangeiro em experimentos deste tipo, e também das facilidades disponíveis para sua execução na instituição de destino, inclusive para condução de experimentos em língua inglesa nativa.

### *6-Avaliação*

Os resultados do experimento serão submetidos à análise estatística de significância para validação da proposta implementada e/ou dos algoritmos utilizados para fins de comparação. Dado o grande volume de dados que se espera coletar, entretanto, uma análise mais completa dos resultados deverá ser realizada por ocasião do retorno ao Brasil.

### *7-Divulgação*

Durante o período de execução do projeto os resultados de cada etapa do trabalho serão divulgados junto ao grupo de pesquisa na instituição de destino e, posteriormente, em eventos científicos da área.

A seguir apresentamos um cronograma resumido para as atividades 1-7, cobrindo o período de 6 meses de execução.

Atividade	Ago	Set	Out	Nov	Dez	Jan
1-Revisão bibliográfica						
2-Refinamento da proposta						
3-Implementação						
4-Definição de experimentos						
5-Execução de experimentos						
6-Avaliação						
7-Divulgação						

### **3.5. Colaboração externa**

A presente proposta prevê que as atividades estipuladas sejam desenvolvidas em colaboração com uma instituição estrangeira de grande tradição em GLN/GER. Pesaram nesta decisão, além da oportunidade de interação com especialistas da

área, o fato de que a condução de experimentos em língua inglesa nativa nos proporciona uma comparação mais direta com o estado da arte, e aumenta drasticamente as possibilidades de divulgação da pesquisa em nível internacional.

A proposta foi assim elaborada com o intuito de ser desenvolvida junto ao grupo de Geração de Língua Natural vinculando ao Departamento de Ciência da Computação da University of Aberdeen<sup>3</sup>, em Aberdeen (Reino Unido), com a colaboração do Dr. Ehud Reiter, um dos pesquisadores mais influentes da área de GLN, e na verdade um de seus precursores. Além de autor do primeiro livro introdutório sobre o assunto (Reiter & Dale, 2000), o Dr. Reiter é também um dos autores do próprio algoritmo Incremental (Dale & Reiter, 1995) que é a base da maioria das abordagens de GER existentes. Sua produção científica abrange praticamente todos os aspectos da pesquisa em GLN, em especial na arquitetura de sistemas deste tipo (Reiter, 2007) e aplicações práticas de GLN de grande porte (Portet et. al., 2009). O Dr. Ehud Reiter foi também o avaliador externo da tese de doutoramento do proponente desta pesquisa (Paraboni, 2003) .

De forma mais diretamente relevante para esta proposta, destacamos que o Dr. Ehud Reiter foi também um dos responsáveis por uma série de estudos sobre o uso de propriedades espaciais em GER (Turner et. al., 2008, 2009) e participa atualmente de um projeto sobre o fenômeno de referência em domínios visuais (Mitchell et. al., 2010) de grande afinidade com nossos atuais interesses. Além disso, seu grupo possui expressiva produção científica na área e tem se destacado em anos recentes por um grande número de pesquisas na área de GER, incluindo o projeto TUNA, que produziu a série de experimentos e o corpus de expressões de referência de mesmo nome (Gatt et. al., 2007) que se tornou um importante *benchmark* da área de GER (Belz et. al., 2007)<sup>4</sup> .

Com relação ao tempo de duração do trabalho, o qual poderia ser considerado curto para projetos desta natureza, destacamos que, em virtude de experiência prévia de doutoramento em outra instituição do Reino Unido (a University of Brighton, de 1998 a 2003), não estamos considerando a necessidade de um período significativo de adaptação, o qual talvez fosse levado em conta no caso de um pesquisador sem experiência prévia no exterior. Além disso, pelo fato de já conhecermos pessoalmente o colaborador principal do projeto e vários membros do

---

<sup>3</sup> <http://www.csd.abdn.ac.uk/research/nlg/>

grupo de GLN da University of Aberdeen, acreditamos que seja perfeitamente viável atingir os objetivos propostos no prazo estipulado.

### **3.6. Resultados Esperados**

O projeto deverá apresentar como produto principal o avanço do estado da arte em GER na forma de novos algoritmos que levem em conta algumas das peculiaridades do fenômeno de referência em domínios físicos mais realistas, como por exemplo a questão da tridimensionalidade do ambiente e o uso de relações espaciais entre objetos. Além disso, espera-se que estreitar os laços de colaboração com o grupo da instituição de destino, neste e em futuros projetos da área, e com isso auxiliar na consolidação da linha de pesquisa em GLN na USP/EACH.

## **4 Considerações Finais**

Este documento apresentou uma proposta de extensão do trabalho prévio deste autor, enfocando a adaptação dos algoritmos de GER apresentados em Paraboni et. al. (2007) e Lucena et. al. (2010) para domínios ditos físicos, como ambientes virtuais tridimensionais e outros. A pesquisa proposta possibilita o avanço do conhecimento na área de GER ao explorar, ao nosso ver de forma inédita, uma abordagem formal de seleção de conteúdo redundante para fins de identificação de referentes neste tipo de contexto, e vem a complementar um projeto atualmente em andamento da instituição de destino que tem como responsável um dos pesquisadores mais influentes da área de GLN/GER.

Este projeto é também parte dos esforços para implantação da linha de pesquisa em GLN na USP/EACH, e complementa nossos estudos em áreas correlatas como planejamento de GLN (Oliveira et. al., 2009), avaliação de sistemas deste tipo (Novais et. al. 2009; Araujo et. al., 2010) e realização superficial estatística (Pereira & Paraboni, 2007,2008; Santos et. al.,2008; Tadeu et. al. 2010; Novais et. al., 2010, 2010a, 2010b).

Destacamos ainda que a instituição de origem passou a contar recentemente (setembro 2010) com um programa de pós-graduação no qual o autor desta proposta tem atuado como docente e orientador de projeto de mestrado na área de

---

<sup>4</sup> Os resultados do projeto TUNA foram tomados como base também para os experimentos apresentados em Lucena et. al. (2010)

concentração da presente proposta<sup>5</sup>. Assim, acreditamos que os resultados deste projeto tragam, além dos avanços científicos esperados, benefícios também para a pesquisa que está sendo desenvolvida localmente e para a internacionalização do novo programa.

## Bibliografia

- Araujo, Roberto Paulo Andrioli de, Rafael Lage de Oliveira, Eder Miranda de Novais, Thiago Dias Tadeu, Daniel Bastos Pereira e Ivandré Paraboni (2010) *SINotas: the Evaluation of a NLG Application*. 7<sup>th</sup> International Conference on Language Resources and Evaluation (LREC-2010) Valletta, Malta, 17-23 maio, pp. 2388-2391.
- Belz, Anja e Albert Gatt (2007) *The Attribute Selection for GRE Challenge: Overview and Evaluation Results*. UCNLG+MT pp. 75-83.
- Byron, D., Alexander Koller, Jon Oberlander, Laura Stoia, and Kristina Striegnitz (2007) *Generating Instructions in Virtual Environments (GIVE): A challenge and evaluation testbed for NLG*. In Proceedings of the Workshop on Shared Tasks and Comparative Evaluation in Natural Language Generation, Arlington.
- Dale, R. e N. Haddock (1991) *Content determination in the generation of referring expressions*. Computational Intelligence, 7(4).
- Dale, R. e E. Reiter (1995) *Computational interpretations of the Gricean maxims in the generation of referring expressions*. Cognitive Science (19).
- Dale, R. e J. Viethen (2009) *Referring Expression Generation through Attribute-Based Heuristics*. 12<sup>th</sup> European Workshop on Natural Language Generation, pp. 58-65.
- Gatt, A., A. Belz e E. Kow (2008) *The TUNA-REG Challenge 2008: Overview and Evaluation Results*. 5<sup>th</sup> International Natural Language Generation Conference (INLG-2008) Salt Fork, USA.
- Gatt, A., A. Belz e E. Kow (2009) *The TUNA-REG Challenge 2009: Overview and Evaluation Results*. 12<sup>th</sup> European Workshop on Natural Language Generation (EACL / ENLG 2009) Atenas.
- Gatt, Albert, I. van der Sluis e Kees van Deemter (2007) *Evaluating algorithms for the generation of referring expressions using a balanced corpus*. Proceedings of the 11<sup>th</sup> European Workshop on Natural Language Generation, pp. 49-56.
- Grice, H. P. (1975) *Logic and Conversation*. In P. Cole and J. L. Morgan (eds.) Syntax and Semantics, Vol. iii: Speech Acts. New York, Academic Press, pp. 41-58.
- Horacek, Helmut (2005) *Generating referential descriptions under conditions of uncertainty*. 10<sup>th</sup> European workshop on Natural Language Generation (ENLG-2005). Aberdeen, 58-67.
- Krahmer, E. e M. Theune (2002) *Efficient Context-Sensitive Generation of Referring Expressions*. In Information Sharing Reference and Presupposition in Language Generation and Interpretation. Kees van Deemter and Rodger Kibble (eds.) CSLI Publications, Stanford, California, pp. 223-264.
- Krahmer, E., S. van Erk e A. Verleg (2003) *Graph-based Generation of Referring Expressions*. Computational Linguistics 29(1).
- Lucena, Diego Jesus de, e Ivandré Paraboni (2008) *Frequency-based Greedy Attribute Selection for Referring Expressions Generation*. 5<sup>th</sup> Intl. Natural Language Generation Conference (INLG-2008). Salt Fork, Ohio.
- Lucena, Diego Jesus de, e Ivandré Paraboni (2008a) *Combining Frequent and Discriminating Attributes in the Generation of Definite Descriptions*. 11<sup>th</sup> Ibero-American Conference on Artificial Intelligence (IBERAMIA-2008) Lisboa. LNAI 5290, pp. 252-261.
- Lucena, Diego Jesus de, e Ivandré Paraboni (2009) *The Design of an Experiment in Anaphora Resolution for Referring Expressions Generation*. Recent Advances in Natural Language Processing (RANLP-2009). Borovets, Bulgaria, pp. 225-229.
- Lucena, Diego Jesus de, e Ivandré Paraboni (2009a) *Improved Frequency-based Greedy Attribute Selection*. 12<sup>th</sup> European Workshop on Natural Language Generation (EACL / ENLG-2009).
- Lucena, Diego Jesus de, e Ivandré Paraboni (2010) *The Anaphor-Antecedent Match: Issues for Referring Expressions Generation*. Proc. of PROPOR-2010.

---

<sup>5</sup> "Geração de Expressões de Referência em Ambientes Virtuais 3D", de Diego dos Santos Silva.

- Lucena, Diego Jesus de, Daniel Bastos Pereira e Ivandré Paraboni (2010) *From Semantic Properties to Surface Text: the Generation of Domain Object Descriptions*. In: *Inteligencia Artificial* 14(45) pp. 48-58.
- Mitchell, M., Kees van Deemter e Ehud Reiter (2010) *Natural Reference to Objects in a Visual Domain*. Proceedings of INLG-2010, Trim, Co. Meath, Irlanda.
- Novais, Eder Miranda de, Rafael Lage de Oliveira, Daniel Bastos Pereira, Thiago Dias Tadeu e Ivandré Paraboni (2009) *A Testbed for Portuguese Natural Language Generation*. 7<sup>th</sup> Brazilian Symposium in Information and Human Language Technology (STIL-2009). São Carlos, 7-11 de setembro.
- Novais, Eder Miranda de, Thiago Dias Tadeu e Ivandré Paraboni (2010) *Text Generation for Brazilian Portuguese: the Surface Realization Task*. NAACL-HLT 2010 Young Investigators Workshop on Computational Approaches to Languages of the Americas . Los Angeles, USA, 1-6 junho, pp.125-131.
- Novais, Eder Miranda de, Thiago Dias Tadeu e Ivandré Paraboni (2010a) *Improved Text Generation using N-gram Statistics*. 12<sup>nd</sup> Ibero-American Conference on Artificial Intelligence (IBERAMIA-2010) Bahía Blanca, Argentina, 1-5 novembro. Springer-Verlag LNAI (no prelo).
- Novais, Eder Miranda de, Thiago Dias Tadeu e Ivandré Paraboni (2010b) *Text-to-Text Surface Realisation using Dependency-Tree Replacement*. 12<sup>nd</sup> Ibero-American Conference on Artificial Intelligence (IBERAMIA-2010) Bahía Blanca, Argentina, 1-5 novembro. Springer-Verlag LNAI (no prelo).
- Oliveira, Rafael Lage de, Eder Miranda de Novais, Roberto Paulo Andrioli de Araujo e Ivandré Paraboni (2009) *A Classification-driven Approach to Document Planning*. Recent Advances in Natural Language Processing (RANLP-2009). Borovets, Bulgaria, 14-16 de setembro, pp. 324-329.
- Paraboni, I. (2003) *Generating References in Hierarchical Domains: the case of Document Deixis*. University of Brighton, tese de doutorado.
- Paraboni, I. e K. van Deemter (2006) *Referring via document parts*. 7<sup>th</sup> International Conference on Intelligent Text Processing and Computational Linguistics (CICLing-2006). Springer Lecture Notes in Computer Science. Cidade do México, pp. 299-310.
- Paraboni, I., J. Masthoff e K. van Deemter (2006) *Overspecified reference in hierarchical Domains: measuring the benefits for readers*. 4<sup>th</sup> Intl. Natural Language Generation Conference (INLG-2006).
- Paraboni, I., K. van Deemter e J. Masthoff (2007) *Generating Referring Expressions: Making Referents Easy to Identify*. Computational Linguistics 33(2) , pp. 229-254.
- Pereira, Daniel Bastos e Ivandré Paraboni (2008) *Statistical Surface Realisation of Portuguese Referring Expressions*. 6<sup>th</sup> International Conference on Natural Language Processing (GoTAL-2008) Gothenburg, Suécia, 25-27 agosto. Lecture Notes in Artificial Intelligence vol. 5221, pp. 383-392. Springer-Verlag.
- Portet, F., E. Reiter, A. Gatt, J. Hunter, S. Sripada, Y. Freer e C. Sykes (2009) *Automatic Generation of Textual Summaries from Neonatal Intensive Care Data*. Artificial Intelligence 173, pp. 789-816.
- Reiter, E. e R. Dale (2000) *Building natural language generation systems*. Cambridge University Press.
- Reiter, Ehud (2007) *An Architecture for Data-to-Text Systems*. Proc. of ENLG-2007, pp. 97-104.
- Santos, Francis Marques Veras dos, Daniel Bastos Pereira e Ivandré Paraboni (2008) *Rule-based vs. Probabilistic Surface Realisation of Definite Descriptions*. VI Workshop on Information and Human Language Technology (TIL-2008). XIV Brazilian Symposium on Multimedia and the Web, pp.372-374.
- Tadeu, Thiago Dias, Eder Miranda de Novais e Ivandré Paraboni (2010) *Extracting Surface Realisation Templates from Corpora*. 7<sup>th</sup> International Conference on Language Resources and Evaluation (LREC-2010) Valletta, Malta, 17-23 maio, pp. 2392-2395.
- Turner, R., S.Sripada, E.Reiter e I. Davy (2008). *Using Spatial Reference Frames to Generate Grounded Textual Summaries of Georeferenced Data*. In Proceedings of INLG08, pp.16-24.
- Turner,R., Y. Sripada, E. Reiter (2009) *Generating Approximate Geographic Descriptions*. Proceedings of ENLG-2009, pp.42-49.
- van Deemter, Kees., I. van der Sluis e A. Gatt (2006) *Building a semantically transparent corpus for the generation of referring expressions*. Proceedings of INLG-2006, Sydney, Australia.