



# Tratamento Computacional da Personalidade Humana para Aplicações de Processamento de Língua Natural

Ivandr  Paraboni (USP / EACH)

**Resumo:** O tratamento computacional de traços de personalidade, seja para reconhecimento destes traços a partir de texto, ou para geração de texto adaptado a um conjunto de traços específico, é um tema central para desenvolvimento de aplicações de Processamento de Língua Natural (PLN) e áreas correlatas. Conhecer os traços de personalidade de um indivíduo (por exemplo, a partir de suas publicações em redes sociais) permite a produção de conteúdo personalizado de várias formas, seja para a apresentação de um website de modo a se tornar mais atraente, para a geração de propaganda mais eficaz e muitas outras. De forma análoga, saber como um indivíduo com determinados traços de personalidade se expressa em língua natural permite a reprodução deste comportamento em aplicações de Geração de Língua Natural (GLN), como a modelagem de personagens realistas de videogames, tutores inteligentes etc. Este documento apresenta uma proposta de projeto de pesquisa na área de PLN enfocando o tratamento de traços de personalidade sob as óticas da interpretação e da geração de língua natural. Tomando por base o modelo dos Cinco Grandes fatores *CGF* amplamente adotado na Psicologia, o projeto prevê a coleta e anotação de um recurso linguístico-computacional básico - que pode ser visto como um corp s paralelo de textos e inventários de personalidade - para mapeamento de relações entre traços de personalidade e fenômenos linguístico variados, e o uso deste recurso para proposta de modelos computacionais de reconhecimento de personalidade a partir de texto, e de geração de texto com base em traços de interesse.

**Palavras-chave:** Processamento de Língua Natural, Traços de Personalidade, Cinco Grandes Fatores.

# 1 *Introdução*

A crescente complexidade dos sistemas computacionais tem sido acompanhada do desenvolvimento de técnicas cada vez mais sofisticadas de comunicação homem-máquina. Sistemas da atualidade são capazes de interpretar e reproduzir uma ampla gama de comportamentos humanos, incluindo, por exemplo, emoções e sentimentos. Estas manifestações de caráter temporário são no entanto decorrentes de um conjunto mais estável e reduzido de características individuais, representando padrões do comportamento humano que são em grande parte previsíveis. Este conjunto de características - ou traços - constitui o que entendemos por *personalidade humana* (Allport e Allport, 1921).

O tratamento computacional da personalidade humana está no centro do projeto de sistemas ditos inteligentes, e será o foco da presente proposta de pesquisa. Traços de personalidade essenciais podem ser determinados com base em diversos métodos propostos pela literatura em Psicologia. Dentre estes, os mais difundidos são os baseados na hipótese lexical, que estabelece que os traços relevantes para caracterização da personalidade são observáveis nas *palavras* que empregamos ao nos comunicar. Esta abordagem foi refinada a partir de um levantamento inicial de 4.500 traços identificados na década de 30 até produzir, de forma independente e simultânea em diversos estudos, um framework estável conhecido como o modelo dos *Cinco Grande Fatores* (*CGF*) ou *Big Five* (Goldberg, 1990).

O modelo *CGF* contempla cinco dimensões fundamentais da personalidade humana - Abertura para experiência, Conscienciosidade, Extroversão, Amabilidade e Neuroticismo - que são amplamente aceitas como a base adequada para a representação da personalidade humana (Andrade, 2008). Do ponto de vista da Ciência da Computação, o modelo *CGF* é de interesse imediato para diversos tipos de estudos de interação humano-computador. Além disso, tendo em vista sua fundamentação linguística, constitui também uma base teórica sólida para uma ampla gama de estudos do Processamento de Língua Natural (PLN). Conforme discutido em (Mairesse et al., 2007), usuários de sistemas computacionais não só atribuem características humanas às máquinas com as quais interagem, como também preferem aquelas que demonstram ter uma personalidade semelhante à sua própria. Esta relação de atração por afinidade sugere assim a necessidade de tratamento computacional da personalidade em duas frentes distintas, porém relacionadas, da pesquisa em PLN: o reconhecimento de personalidade a partir de texto, e a geração de texto baseada em traços de personalidade.

O tratamento computacional de traços de personalidade é um tema central para desen-

volvimento de aplicações de PLN. Conhecer os traços de personalidade de um indivíduo (por exemplo, a partir de suas publicações em redes sociais) permite produzir conteúdo personalizado de várias formas, seja para apresentação de um website de modo a se tornar mais atraente, para a geração de propaganda mais eficaz etc. De forma análoga, saber como um indivíduo com certos traços de personalidade se expressa em língua natural permite a reprodução deste comportamento em aplicações de Geração de Língua Natural (GLN). Estas aplicações incluem, por exemplo, a modelagem de personagens realistas de videogames ou filmes, sistemas geradores de narrativas, tutores inteligentes e outras.

Interpretação e geração de língua natural baseadas em traços de personalidade constituem desafios de pesquisa consideráveis e, apesar de sua natureza complementar e por vezes concorrente (por exemplo, em aplicações de diálogo homem-máquina), poderiam em princípio ser tratadas como questões de pesquisa independentes. Sob o ponto de vista prático, entretanto, observamos que uma abordagem como a que será proposta - baseada em métodos de aprendizagem de máquina - parte em ambos os casos de um mesmo pressuposto básico: a construção de mapeamentos entre traços de personalidade e fenômenos linguísticos. Esta base comum sugere um projeto de caráter exploratório para construção de recursos linguístico-computacionais deste tipo, e o uso destes recursos em um estudo de ampla cobertura da relação entre língua natural e traços de personalidade para desenvolvimento de modelos computacionais de interpretação e geração de texto. Uma proposta de pesquisa desta natureza é discutida a seguir.

## 1.1 Objetivos

O objetivo da pesquisa proposta é a construção de um *córpus paralelo* composto de documentos textuais e dos respectivos inventários de personalidade de seus autores, e uso deste *córpus* no desenvolvimento de modelos computacionais de *interpretação* e *geração* de língua natural baseados em traços de personalidade.

Do ponto de vista da *interpretação* de língua natural, a pesquisa contempla o estudo e desenvolvimento de modelos computacionais para reconhecimento de traços de personalidade e outras formas de caracterização autoral a partir de texto produzido em formato livre ou controlado em Português brasileiro.

Do ponto de vista da *geração* de língua natural, a pesquisa contempla o estudo e desenvolvimento de modelos computacionais de produção de texto com base em traços de personalidade sob três perspectivas: a seleção de conteúdo referencial (Krahmer e Deemter, 2012), o planejamento sentencial (Stent e Molina, 2009), e a realização superficial sentencial do tipo texto-para-texto (Reiter e Dale, 2000).

## 1.2 Contribuições previstas

As diversas questões teóricas suscitadas por um estudo do tipo proposto (a serem discutidas na seção a seguir) são atuais e relevantes para a pesquisa em PLN/GLN, e sua investigação representa uma contribuição genuína para o estado da arte nestas áreas. Em especial, acreditamos que o estudo proposto seja o primeiro do gênero a contemplar de forma concomitante a relação entre traços de personalidade e as questões de interpretação e geração de língua natural para o Português do Brasil.

De forma mais específica, o projeto auxiliará também na consolidação da linha de pesquisa em PLN/GLN desta instituição em complemento a estudos como (Paraboni et al., 2007; Lucena et al., 2010; Novais e Paraboni, 2012; Paraboni e Deemter, 2014; Silva e Paraboni, 2015). Além disso, à luz do recente programa de pós-graduação da instituição, espera-se formar pela primeira vez uma equipe de estudantes de pós-graduação atuando em projetos complementares vinculados a esta linha de pesquisa, assim como estreitar os laços de colaboração com instituições com a qual compartilhamos interesses neste tema.

## 2 *O tratamento computacional da personalidade humana*

Nesta seção apresentamos uma visão geral do tratamento computacional de personalidade sob a ótica da interpretação e geração de língua natural, e alguns de seus desafios em aberto.

### 2.1 Aquisição de conhecimento

O tratamento computacional de personalidade a partir de texto coloca de imediato a questão da aquisição do conhecimento necessário à tarefa, ou seja, de como obter material textual anotado com informações de personalidade. Enquanto aplicações mais tradicionais como a classificação de documentos podem se beneficiar de bases textuais previamente rotuladas, a base textual aqui exigida não possui equivalente disponível para o caso do Português do Brasil. A necessidade de um recurso deste tipo - que pode ser visto como um *córpus* paralelo contendo (a) informações de personalidade e (b) textos - é discutida a seguir.

Com relação às informações de personalidade (a), consideramos nesta proposta o modelo dos *Cinco Grande Fatores (CGF)* ou *Big Five* (Goldberg, 1990). O modelo *CGF* contempla cinco dimensões fundamentais da personalidade humana - Abertura para experiência, Conscientiosidade, Extroversão, Amabilidade e Neuroticismo - que são amplamente aceitas como a base adequada para a representação da personalidade humana segundo pelo menos quatro perspectivas (Andrade, 2008): (1) estudos longitudinais e de observação cruzada têm demonstrado que os cinco fatores são disposições duradouras manifestas em padrões de comportamento; (2) os traços relatados por cada fator são encontrados em uma variedade de teorias de personalidade, bem como na linguagem usual de descrição destes; (3) os fatores são encontrados em diferentes idades, sexos, raças e nacionalidades, embora variem em certo grau segundo a cultura; (4) evidências da hereditariedade sugerem que os cinco fatores têm uma base biológica que não seria diretamente influenciada pelo ambiente.

Feitas estas observações, consideramos que o *córpus* paralelo exigido para o estudo proposto deva contemplar a anotação de uma base textual com os cinco fatores previstos no modelo *CGF*. Estes fatores podem ser estimados por diversos métodos consagrados em Psicologia, sendo o mais comum o uso de inventários de traços de personalidade. Inventários deste tipo são questionários de personalidade aplicados por um profissional da área, ou respondidos de forma autônoma pelo próprio sujeito avaliado. Um exemplo proeminente de inventário para captura dos fatores

de personalidade *CGF* é o *NEO-PI-R*, que consiste de 240 itens (perguntas), e possui uma versão reduzida *NEO-FFI* de 60 itens. Ambos inventários permitem a diferenciação de cada um dos cinco fatores fundamentais em termos de facetas mais específicas. Por exemplo, o fator Extroversão pode ser subdividido em facetas como ‘sociabilidade’, ‘assertividade’ e outras.

A necessidade de um instrumento de avaliação mais ágil, entretanto, levou à proposta do inventário *BFI* (John et al., 1991). Este inventário consiste de 44 itens na forma de frases breves contendo adjetivos que capturam os aspectos mais essenciais de cada fator do modelo *CGF*. A versão em Inglês do inventário *BFI* foi desenvolvida a partir de análise fatorial de grandes massas de dados e, apesar do número reduzido de itens, é considerada uma medida segura dos atributos mais importantes do modelo *CGF* sem prejuízo à cobertura ou às boas qualidades psicométricas observadas em inventários mais extensos, sendo pelo menos tão eficiente e fácil de interpretar quanto o inventário *NEO-FFI*. O inventário *BFI* é considerado especialmente atraente para aplicações computacionais, e testes mais extensos são geralmente recomendados apenas para os casos em que a disponibilidade de tempo dos participantes não é um obstáculo, quando estes participantes possuem boa escolaridade e experiência em testes do gênero, e quando a pesquisa requer o exame de múltiplas facetas do modelo (John et al., 2008).

O inventário *BFI* tem sido replicado em dezenas de outros idiomas, incluindo alguns estudos dedicados ao Português. Em especial, o estudo em (Andrade, 2008) validou o *BFI* para o Português brasileiro por meio de uma análise fatorial envolvendo uma amostra de 5.089 respondentes das cinco regiões brasileiras. Não pudemos identificar dentre estes estudos, entretanto, nenhum caso em que os inventários fossem acompanhados de material textual na quantidade e teor exigidos para a presente pesquisa - o que é natural considerando-se que o foco destes estudos não costuma ser o reconhecimento de personalidade a partir de texto - e constatamos assim que um córpus paralelo do tipo pretendido não se encontra disponível para pesquisa em PLN do Português brasileiro<sup>1</sup>.

Deixando-se de lado a questão do inventário de personalidade, resta ainda a questão (b) de qual o tipo de texto que o córpus pretendido deveria contemplar, e como de fato obtê-lo. Em estudos como (Mairesse et al., 2007), por exemplo, múltiplas fontes textuais são consideradas, incluindo relatos de fluxo de consciência - onde indivíduos escreviam por 20 minutos a respeito de qualquer coisa que lhes ocorresse em mente - extraídos de um conjunto de 2400 redações anotadas com as dimensões *CGF*.

Outras fontes textuais comuns para computação de personalidade são os blogs pessoais (Oberlander e Nowson, 2006; Nowson e Oberlander, 2007; Yarkoni, 2010; Iacobelli et al., 2011) e redes sociais como Twitter (Qiu et al., 2012; Nunes et al., 2013) ou Facebook (Schwartz et al.,

---

<sup>1</sup>Uma exceção é o trabalho em (Nunes et al., 2013), no qual foram coletadas publicações na rede social Twitter acompanhadas de um inventário de personalidade de 28 participantes. Este volume de dados entretanto não viabilizaria um projeto como o do tipo aqui pretendido.

2013). Para a plataforma Facebook há inclusive uma base dedicada à computação de personalidade - denominada *myPersonality* (Kosinski et al., 2015) - composta de uma vasta coleção de atualizações de status anotada com os cinco fatores e diversos outros tipos de informações. Esta base conta atualmente com mais de 4 milhões de perfis de usuários de língua inglesa.

Textos provenientes de redes sociais são de fácil acesso e constituem uma fonte potencialmente útil para diversos tipos de aplicações computacionais, e serão por este motivo considerados na proposta discutida na Seção 3.1.1. Conforme observado em (Celli, 2012), entretanto, diferentes domínios textuais tornam explícitos diferentes aspectos da personalidade, e considerando-se os interesses específicos e o histórico de pesquisa em GLN do grupo responsável por esta proposta, observamos que este tipo de texto pode não ser suficiente para um estudo mais aprofundado de certas questões de pesquisa da área.

De forma mais específica, observamos que textos produzidos de forma livre (seja na forma de redações ou publicações em redes sociais) podem não fornecer subsídio para o estudo de questões de GLN como as que serão discutidas na Seção 2.3 (a saber, a seleção de conteúdo referencial, o planejamento sentencial e a realização superficial na geração texto-para-texto), e que são de especial interesse para a presente proposta. Para estudos deste tipo, assim como em vários outros problemas de GLN, coloca-se a necessidade da observação da produção de língua humana sob condições *controladas*, de modo que seja possível observar não apenas o resultado final (i.e., a língua produzida), mas também as *condições iniciais* (ou estímulos) que a motivaram.

Com base nestas observações, consideramos assim a necessidade de construção de uma base de textos em Português com informações de personalidade que seja de propósito geral, e que ofereça suporte para vários tipos de pesquisa em PLN e GLN. Uma base deste tipo deveria idealmente contemplar não apenas textos redigidos em formato livre e não-estruturado (como os provenientes de redes sociais), mas também alternativas em que o pesquisador possa exercer maior controle sobre o estímulo inicial, obtendo assim exemplos de produção linguística mais estruturada.

## 2.2 O reconhecimento de personalidade a partir de texto

Supondo-se que uma base de conhecimento como a discutida na seção anterior esteja de alguma forma disponível, esta seção discute alguns desafios relacionados ao reconhecimento automático de personalidade. Ao contrário de uma tarefa tradicional de classificação de documentos, a tarefa computacional de reconhecimento de personalidade é menos baseada no conteúdo do texto, e mais baseada na sua forma (representada, por exemplo, pela sua variação estilística). Estudos deste tipo costumam seguir uma metodologia tradicional de aprendizagem de máquina supervisionada (Oberlander e Nowson, 2006) ou semi-supervisionada (Celli, 2012) para reconhecimento das dimensões *CGF* ou de outras características autorais. Este reconhecimento pode

assumir a forma de um problema de classificação (e.g., binária, decidindo se um indivíduo é extrovertido ou não), de um problema de regressão (e.g., determinando o valor escalar da dimensão Extroversão do indivíduo) ou ainda de um problema de *ranking* (e.g., ordenando um conjunto de indivíduos segundo uma dimensão de interesse). De modo geral, entretanto, o problema tem sido frequentemente abordado de forma limitada (e.g., considerando apenas um ou dois fatores do modelo *CGF*), e mesmo assim com resultados relativamente modestos.

Uma das primeiras iniciativas de grande escala para tratamento computacional do reconhecimento de personalidade a partir de texto é o trabalho em (Argamon et al., 2005), que conduziu um experimento envolvendo 2263 ensaios escritos por 1200 estudantes que haviam preenchido o questionário *NEO-FFI*, porém limitado aos extremos inferior e superior da escala para Extroversão e Neuroticismo. As palavras dos ensaios foram agrupadas em quatro categorias de significado psicológico definido: funções (artigos, preposições etc.), coesão (demonstrativos etc.), avaliação (termos que avaliam o conteúdo quanto à validade, verossimilhança, aceitação etc.) e julgamento (termos que expressam a atitude do autor em relação ao conteúdo). Os textos foram representados pelas frequências relativas de cada categoria, e as classes binárias Extroversão e Neuroticismo foram classificadas usando SVMs, com uma acurácia máxima de 58%.

Em (Mairesse et al., 2007), uma versão expandida do mesmo conjunto de textos e inventários utilizado em (Argamon et al., 2005) foi empregada em uma abordagem que explora 88 categorias de palavras extraídas da base psicolinguística LIWC (*Linguistic Inquiry and Word Count*) (Tausczik e Pennebaker, 2010) e 26 atributos da base MRC (*Medical Research Council*) composta de 150.837 itens lexicais. Estes atributos incluíram várias normas como concretude, idade de aquisição de palavras e outras. O experimento realizado consistia em discriminar os extremos superior e inferior para as cinco dimensões *CGF*, obtendo acurácia máxima de 50% a 62% com uma abordagem baseada em SVMs.

Estudos como em (Argamon et al., 2005; Mairesse et al., 2007) adotam uma abordagem lexical baseada em estatísticas sobre o uso de palavras individuais. Estudos como (Oberlander e Nowson, 2006; Nowson e Oberlander, 2007), por outro lado, dispensam este tipo de conhecimento fazendo uso de modelos de n-gramas. Nestes estudos, o objetivo foi mais uma vez o de discriminar indivíduos de pontuação alta e baixa para quatro dos cinco fatores de personalidade (excetuando-se Abertura). A classificação fez uso do algoritmo Naive-Bayes e de SVMs. Em (Oberlander e Nowson, 2006) foi utilizado um conjunto de 71 blogs, obtendo-se acurácia de 45% (aleatório) até 100% dependendo do modelo de n-gramas utilizado e das classes definidas. Em (Nowson e Oberlander, 2007), o mesmo experimento foi repetido utilizando-se um conjunto de 1672 blogs, com acurácia máxima de 65%.

O reconhecimento de traços fundamentais de personalidade possui um certo grau de afinidade com o problema correlato de caracterização autoral, aqui entendido como o reconhecimento



de outros traços identificadores dos autores do texto além do modelo *CGF*, e sua possível relação com este. Exemplos de caracterização autoral desta natureza incluem estimativas de idade e gênero (Marquardt et al., 2014), afiliação política e agressividade verbal (Bates et al., 2012) etc.

O problema computacional do reconhecimento de personalidade a partir de texto é bastante desenvolvido para o caso do idioma Inglês. Usando dados como os provenientes de bases públicas como *myPersonality* (Kosinski et al., 2015), a área já conta inclusive com alguns eventos científicos dedicados, como a série de ‘shared tasks’ PAN<sup>2</sup>. Nenhuma destas iniciativas, entretanto, encontra similar para o Português brasileiro no contexto do modelo *CGF*. Além da escassez de dados discutida na seção anterior, observa-se também a inexistência de sistemas ou modelos que sirvam de termo de comparação (ou *baseline*) para uma pesquisa desta natureza. Diante deste cenário, sugere-se a necessidade de um estudo exploratório do gênero, no qual sejam organizados os principais recursos necessários, e seja realizado um mapeamento entre traços de personalidade e fenômenos linguísticos de interesse para o PLN do Português brasileiro.

## 2.3 Geração de texto baseada em traços de personalidade

Esta seção discute o problema ‘inverso’ ao reconhecimento de personalidade tratado na seção anterior, ou seja, a questão de como expressar traços de personalidade na produção automática de texto. Sistemas deste tipo - denominados sistemas de Geração de Língua Natural (GLN) (Reiter e Dale, 2000) - produzem descrições textuais a partir de uma entrada de dados geralmente não linguística, e são empregados quando o uso de texto predefinido não é suficiente, ou seja, quando é necessária uma maior variação linguística nos documentos gerados e/ou maior proximidade em relação ao desempenho humano.

Diversos estudos psicolinguísticos abordam a relação entre personalidade e produção da língua. Do ponto de vista computacional, entretanto, é observado em (Mairesse e Walker, 2011) que a maioria dos sistemas de GLN existentes trata prioritariamente ou exclusivamente da tarefa de produzir uma saída única e gramaticalmente correta que satisfaça os objetivos da comunicação. Nos casos em que variação estilística é considerada (Paiva e Evans, 2004; Reiter e Williams, 2010), normalmente não é considerado o objetivo de modelar um locutor específico ou, como no caso da presente proposta, um perfil psicológico específico.

Dentre os poucos estudos de GLN baseada em traços de personalidade, destaca-se o sistema PERSONAGE em (Mairesse e Walker, 2010, 2011). Neste estudo é apresentado um método parametrizável motivado por correlações observadas entre traços de personalidade do modelo *CGF* e uma ampla gama de decisões de geração. O sistema é treinado a partir de um corpus paralelo de textos e inventários de personalidade semelhante ao discutido na Seção 2.1, e é capaz

---

<sup>2</sup><http://pan.webis.de/tasks.html>

de expressar uma série de diferenças de personalidade que são perceptíveis por leitores humanos.

Apesar do grau de sofisticação da abordagem proposta em (Mairesse e Walker, 2011), entretanto, observamos que o sistema faz uso de diversos componentes pré-existentes da arquitetura GLN específicos para o Inglês, o que inviabiliza o uso destas mesmas técnicas no presente estudo. Além disso, sistemas como PERSONAGE seguem uma arquitetura tradicional de GLN do tipo dados-para-texto (Reiter e Dale, 2000), ou seja, partem de uma representação não-linguística como entrada para então construir o texto de saída. As vantagens deste tipo de arquitetura para diversos tipos de aplicações são bem conhecidas, mas suas limitações também. Dentre elas, a mais importante é a questão do formato dos dados de entrada, já que sistemas deste tipo são tipicamente dependentes do domínio da aplicação e, conseqüentemente, pouco reutilizáveis.

Neste cenário, diversos estudos recentes têm privilegiado o desenvolvimento de aplicações de geração texto-para-texto nos quais a própria entrada é representada em língua natural, e a tarefa de geração pode assim ser vista como um processo de reescrita. No entanto, não encontramos na literatura estudos sobre sistemas de geração texto-para-texto que considerem a questão da personalidade humana. Sistemas de reescrita deste tipo poderiam, por exemplo, impor um perfil psicológico de interesse a um texto previamente escrito por humanos, constituindo uma provável solução de baixo custo para a modelagem de personagens realistas de videogame, ou para a personalização de roteiros de diálogos homem-máquina em geral.

Com base nestas observações, destacamos - de forma obviamente não-exaustiva - três lacunas de pesquisa de especial interesse para a presente proposta, e que não são abordadas no trabalho prévio de GLN baseada em traços de personalidade: (a) a seleção de conteúdo referencial, (b) o planejamento sentencial e (c) a realização superficial do tipo texto-para-texto. Estas questões são discutidas brevemente a seguir, e de forma mais detalhada na Seção 3.1.

A seleção de conteúdo referencial (a) é a tarefa de determinar quais propriedades semânticas serão expressas na forma de descrições definidas ou indefinidas como em ‘a mulher de preto’, ‘Dona Maria do terceiro andar’, ‘uma senhora de óculos, com um cachorro’. A geração computacional e psicologicamente plausível de expressões deste tipo é um tópico recorrente na pesquisa em GLN (Krahmer e Deemter, 2012), e tem sido regularmente abordado pelo autor da presente proposta (Paraboni et al., 2007; Paraboni e Deemter, 2014; Silva e Paraboni, 2015; Paraboni et al., 2016). Embora seja reconhecido que, assim como em inúmeras outras tarefas de produção de língua natural, a estratégia de seleção de conteúdo varia de pessoa para pessoa (Ferreira e Paraboni, 2014), não há, até onde temos conhecimento, nenhum estudo sobre como diferentes tipos de personalidade desempenham esta tarefa.

A geração texto-para-texto contempla também uma série de subtarefas relacionadas ao planejamento sentencial (b) (Stent e Molina, 2009), como questões de ordenação sentencial, agregação e inserção de marcadores de discurso e relacionadas ao próprio planejamento discursivo.

sivo, além da própria realização da forma superficial (c) (Novais e Paraboni, 2012). Técnicas deste tipo podem ser implementadas, por exemplo, com o uso de regras codificadas de forma manual ou extraídas de um *cópus* anotado e têm, em alguns casos isolados, considerado a modelagem do indivíduo (Walker et al., 2007). Novamente, entretanto, não identificamos nenhum estudo de planejamento ou realização sentencial para geração texto-para-texto que leve em conta modelos de personalidade humana.

## 2.4 Considerações

Podemos resumir os desafios do tratamento computacional da personalidade humana discutidos nas seções anteriores em três grandes questões:

- Ausência de recursos linguístico-computacionais na forma de um *cópus* paralelo de inventários de personalidade e textos produzidos em forma livre e controlada.
- Escassez de conhecimento sobre a relação entre personalidade e fenômenos linguísticos, seja para reconhecimento de personalidade a partir de texto, ou para geração de texto com base em traços de personalidade.
- Ausência de estudos sobre a geração texto-para-texto baseada em traços de personalidade, incluindo tarefas específicas como a seleção de conteúdo referencial, o planejamento sentencial e a realização superficial baseada nestes traços.

## 3 *Proposta de pesquisa*

Conforme discutido na Seção 1.1, a pesquisa proposta objetiva a construção de um *córpus paralelo* de documentos textuais e respectivos inventários de personalidade de seus autores, e uso deste *córpus* no desenvolvimento de modelos computacionais de *interpretação* e *geração* de língua natural baseados em traços de personalidade. Esta proposta é detalhada a seguir.

### 3.1 Atividades

O projeto consiste de uma etapa inicial de construção do *córpus*, e duas etapas independentes de investigação do tratamento computacional da personalidade: uma tratando do reconhecimento de traços de personalidade e outras formas de caracterização autoral a partir de texto, e a outra tratando da geração de língua natural com base em traços de personalidade.

#### 3.1.1 Construção do *córpus paralelo*

A primeira etapa do projeto - construção do *córpus paralelo* - consiste da coleta e pré-processamento de inventários de personalidade com base no modelo *CGF* (John et al., 1991, 2008), e de textos produzidos em formato livre e controlado. Esta etapa contempla cinco atividades principais: (1a) o desenvolvimento de ferramentas de coleta de dados; (1b) a aplicação do inventário de personalidade a um grupo de participantes; (1c) a coleta de textos produzidos por estes mesmos participantes de forma livre (i.e., a partir da plataforma Facebook); (1d) a coleta de textos de forma controlada (como resposta a estímulos de um experimento presencial); e (1e) o pré-processamento destes textos para uso nas etapas subsequentes do projeto.

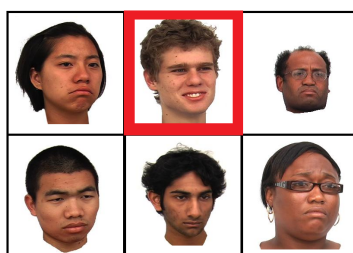
A atividade de desenvolvimento de ferramentas (1a) envolve a criação de um aplicativo para a plataforma Facebook que, além de permitir a resposta ao inventário de personalidade, faz a coleta simultânea das publicações de cada usuário (mediante autorização prévia). Este aplicativo, que segue uma metodologia semelhante à discutida em (Schwartz et al., 2013), será utilizado tanto na aplicação do inventário de personalidade (1b) como na coleta de textos de produção livre (atividade 1c). Para participantes que não sejam usuários Facebook, será disponibilizada também uma versão off-line do inventário, embora obviamente neste caso não haverá coleta de textos provenientes da rede social. Ao invés disso, o inventário será seguido do experimento de coleta de textos de modo controlado (atividade 1d).

O inventário de personalidade a ser utilizado é o *ICGF*, que foi validado para o Português

brasileiro em (Andrade, 2008). Assim como em vários dos trabalhos mais influentes desta área como (Argamon et al., 2005; Oberlander e Nowson, 2006; Mairesse et al., 2007), neste projeto será utilizado o método de auto avaliação da personalidade, ou seja, utilizando-se questionários respondidos pelo próprio sujeito avaliado. Embora estudos como (Mairesse et al., 2007) indiquem que resultados mais precisos podem ser obtidos com o emprego direto de especialistas humanos (i.e., psicólogos) na tarefa, o custo de uma abordagem deste tipo seria excessivo e possivelmente não justificável face ao caráter exploratório do projeto. Além disso, se resultados satisfatórios puderem ser obtidos a partir de dados de auto avaliação, é razoável supor que estes resultados sejam ainda melhores se/quando os modelos propostos puderem ser recriados com dados de avaliações produzidas por especialistas.

Com relação à base textual a ser coletada, assim como em (Schwartz et al., 2013), consideramos o uso de texto não-estruturado proveniente da rede social Facebook. No entanto, conforme discutido na Seção 2.1, observamos que no caso específico da pesquisa em GLN é necessário conhecer não apenas exemplos de produção da língua, mas também os estímulos iniciais que os motivaram. Por este motivo, a base textual a ser construída contemplará também três tipos de texto produzido sob condições controladas. Estes textos serão produzidos por participantes de um experimento presencial em resposta a determinados estímulos visuais de interesse, e estão relacionados a três tarefas de produção de língua natural: (a) a identificação de entidades visuais, (b) a descrição livre e multi-sentencial de imagens; e (c) a produção de legendas descritivas na forma mono-sentencial das mesmas imagens.

A subtarefa de identificação de entidades visuais (a) está diretamente ligada ao problema de seleção de conteúdo referencial de GLN. Experimentos prévios deste tipo realizados nesta instituição levaram, por exemplo, à construção de córpus de descrições definidas como Stars2 (Paraboni et al., 2016). Diferentemente de trabalhos prévios baseados em domínios simplificados (e.g., objetos geométricos), entretanto, o domínio considerado neste caso fará uso de imagens de estímulo com maior potencial de explicitar diferenças entre traços de personalidade. De forma mais específica, os contextos de referência a serem considerados farão uso de imagens extraídas da base *Face Place* (Righi et al., 2012) de fotografias humanas validadas para diversos tipos de emoções e características físicas. Um exemplo de imagem de estímulo a ser utilizada nesta atividade é ilustrado na Fig. 1.



**Figura 1** – Imagem de estímulo criada a partir da base *Face Place* (Righi et al., 2012).

Nesta tarefa, o participante é instruído a referenciar de forma única (i.e., sem ambiguidade) a pessoa ou entidade destacada de tal modo que outra pessoa possa identificá-la. No caso do presente exemplo, isso poderia ser feito, por exemplo, com uso de descrições como ‘o rapaz loiro’ ou ‘o homem que está sorrindo, no centro superior da imagem’. O objetivo deste tipo de coleta de dados é modelar as decisões sobre o conteúdo - ou ‘o que dizer’, cf. (Krahmer e Deemter, 2012) - de expressões produzidas por indivíduos com diferentes traços de personalidade.

O experimento contempla também duas subtarefas de descrição de imagens: (b) em versão detalhada (em texto livre e multi-sentencial) e (c) em versão resumida (na forma de uma sentença única). Os estímulos visuais neste caso serão provenientes da base *GAPED* (Dan-Glauser e Scherer, 2011) de imagens classificadas por valência e significância normativa designadas de modo a despertar diferentes graus de emoção positiva e negativa. Um exemplo de imagem disponibilizada pela base *GAPED* e ilustrado na Fig. 2.



**Figura 2** – Imagem de estímulo da base *GAPED* (Dan-Glauser e Scherer, 2011).

Ao contrário da tarefa de identificação anterior, o objetivo neste caso não é observar a estratégia referencial do sujeito do experimento, mas sim a estratégia utilizada para determinar os elementos mais importantes da imagem, a ordem e estruturação da descrição, e suas escolhas lexicais e sintáticas. Esta coleta de dados será realizada em duas versões - detalhada e resumida - como forma de exercer um maior grau de controle sobre o texto produzido, sem no entanto influenciá-lo. Além disso, as duas versões do texto permitirão o estudo de questões ligadas à produção textual tanto em nível sentencial quanto discursivo, a saber: as descrições detalhadas (b) serão explorada em um estudo sobre o planejamento sentencial (Stent e Molina, 2009) enquanto que as descrições resumidas (c) serão exploradas em um estudo de realização superficial sentencial para geração do tipo texto-para-texto, em ambos os casos levando em conta os traços de personalidade do autor do texto.

No momento da elaboração da presente proposta, protótipos de ambas ferramentas previstas na atividade 1a encontram-se em fase de teste. O tempo estimado para resposta ao inventário de personalidade e coleta simultânea de publicações Facebook é de cerca de 10 minutos. O experimento presencial toma cerca de 40-60 minutos para ser completado, descontando-se o tempo do inventário. As permissões necessárias para realização dos dois tipos de coleta de dados já foram concedidas pelo Comitê de Ética em Pesquisa desta instituição, e são acompanhadas

de um termo de concordância que deve ser aceito pelos sujeitos antes da participação, tanto em modo off-line como via Facebook.

A aplicação do inventário de personalidade (atividade 1b) será realizada preferencialmente via Facebook de modo a aproveitar tanto quanto possível o recurso de coleta automática das publicações do participante (atividade 1c). Estabelecemos a meta de obter um mínimo de 1000 inventários acompanhados de publicações Facebook. Além disso, estabelecemos também a meta de coletar um mínimo de 120 conjuntos de textos produzidos sob condições controladas no experimento presencial (atividade 1d).

Finalmente, a atividade (1e) contempla o pré-processamento dos textos coletados. Em um primeiro momento, será realizada a correção ortográfica e detecção de idioma com apoio de dicionários do Português e Inglês (Muniz, 2004), o tratamento de falsos homógrafos e abreviaturas comuns em textos provenientes de redes sociais (Duran e Nunes, 2015), seguidas de análise de *part-of-speech* (Fonseca e Rosa, 2013) e sintática (Bick, 2000).

O produto final desta etapa é o *cópus* paralelo de textos (livres e controlados) e inventários *ICGF*, enriquecido com informações morfossintáticas e semânticas. Este *cópus* será empregado nas duas etapas subsequentes do projeto. Porções não-confidenciais do *cópus* serão disponibilizadas para reuso pela comunidade científica após a divulgação dos resultados finais do projeto.

### 3.1.2 Reconhecimento de personalidade a partir de texto

A segunda etapa do projeto - o reconhecimento de personalidade a partir de texto - consiste do uso do *cópus* paralelo de textos e inventários de personalidade coletado na etapa anterior em duas atividades de investigação relacionadas à interpretação de língua natural: (2a) o reconhecimento de traços de personalidade básicos a partir de texto de produção livre; e (2b) um estudo sobre outras formas de caracterização autoral a partir destes dados. Estas atividades são complementadas com uma fase de refinamento (2c) dos modelos propostos, e de sua consolidação na forma de uma proposta mais geral para solução destes problemas.

Assim como em (Argamon et al., 2005; Celli, 2012), no presente projeto o reconhecimento de personalidade a partir de texto (2a) será desenvolvido com o uso de técnicas de aprendizagem de máquina supervisionada para classificação e/ou regressão dos cinco fatores do modelo *CGF* a partir dos texto coletado via Facebook. Para este fim, será considerada uma abordagem lexical semelhante à adotada em (Argamon et al., 2005; Mairesse et al., 2007). Esta abordagem será baseada em atributos extraídos do próprio texto ou de bases auxiliares que estejam disponíveis para o Português, como a versão brasileira do dicionário LIWC (Tausczik e Pennebaker, 2010) discutida em (Filho et al., 2013), e já utilizado em estudos como (Nunes et al., 2013).

Outras possíveis fontes de conhecimento a serem consideradas incluem normas de concretude

(Janczura et al., 2007; Calais et al., 2012), alerta e valência (Oliveira et al., 2013) e outras formas de medida lexical (Scarton, 2009). Além destas, consideramos ainda a possibilidade de emprego de modelos baseados em n-gramas como em (Oberlander e Nowson, 2006; Nowson e Oberlander, 2007) que poderiam inclusive vir a ser combinados com o método lexical, sugerindo-se assim uma ampla gama de novas abordagens híbridas para o problema.

A caracterização complementar a que se refere a atividade 2b trata do reconhecimento de outros traços identificadores dos autores do texto além do modelo *CGF*, e sua possível relação com este. Com base nos dados coletados nas etapas anteriores do projeto, consideramos em um primeiro momento a caracterização de traços disponibilizados pela rede social, tais como faixa etária, gênero, formação acadêmico e afins, utilizando métodos de aprendizagem similares aos considerados na atividade 2a.

Finalmente, a atividade 2c tratará do refinamento destas iniciativas, e da sua consolidação na forma de novos modelos computacionais de reconhecimento de traços de personalidade e de outras formas de caracterização autoral a partir de texto.

### 3.1.3 Geração de texto baseada em traços de personalidade

A terceira e última etapa do projeto trata do uso do conhecimento produzido nas etapas anteriores do projeto para proposta de modelos de geração de língua natural texto-para-texto baseada em traços de personalidade, questão de especial afinidade com a pesquisa prévia desenvolvida pelo grupo proponente. Para este fim, serão contempladas três questões de especial interesse, a saber: (3a) a seleção de conteúdo referencial, o planejamento sentencial (3b) e a realização superficial em geração texto-para-texto (3c). Estas atividades serão também seguidas de refinamentos e consolidação (3d).

A atividade 3a consiste em uma investigação de como diferentes tipos de personalidade desempenham a tarefa de seleção de atributos de expressões de referência. O conhecimento proveniente deste estudo será empregado no projeto de um algoritmo de geração destas expressões que considere, além dos parâmetros de entrada usuais (leia-se: um contexto contendo um alvo que se deseja distinguir de outros objetos por meio de suas propriedades semânticas), parâmetros adicionais especificando o tipo de personalidade desejado. Um estudo deste tipo é uma extensão natural de pesquisas recentes realizadas nesta instituição, como (Paraboni e Deemter, 2014; Silva e Paraboni, 2015; Paraboni et al., 2016) e, de forma mais específica, do estudo em (Ferreira e Paraboni, 2014), que tratou da questão da variação humana nesta tarefa.

A atividade 3b consiste em um estudo de como diferentes tipos de personalidade definem a estrutura de um discurso multi-sentencial. Este estudo objetiva a proposta de modelos de planejamento sentencial contemplando questões de ordenação das sentenças e entidades discursivas.





### 3.3 Resultados esperados

O projeto prevê a disponibilização de três produtos principais, a saber:

- (1) Um cópús paralelo de textos e inventário de personalidade de seus autores.
- (2) Um modelo de reconhecimento de traços de personalidade e outras formas de caracterização autoral a partir de textos em Português brasileiro.
- (3) Um modelo de GLN contemplando as questões de seleção de conteúdo referencial, planejamento sentencial e realização superficial com base em traços de personalidade.

O cópús a que se refere o item (1) acima é um recurso linguístico-computacional inédito para o tratamento automático de personalidade em Português brasileiro, e um requisito para a presente e também para futuras pesquisas na área. O cópús consiste dos cinco subconjuntos de dados detalhados na seção 3.1.1, a saber: uma base de inventários de personalidade e quatro tipos de bases textuais, sendo um subcópús de textos em formato livre proveniente de redes sociais, e três subcópús de texto controlado (expressões de referência em um domínio de teor afetivo, e descrições de imagens em formato mono e multi-sentencial).

Os modelos computacionais resultantes deste estudo deverão avançar diversas questões de pesquisa relacionadas ao tratamento automático de traços de personalidade, tanto do ponto de vista da interpretação (item 2 acima) quanto da geração (item 3) de língua natural, conforme discutido nas Seções 3.1.2 e 3.1.3. Estas questões são em grande parte extensão de pesquisas já realizadas nesta instituição, e devem auxiliar o futuro desenvolvimento de aplicações de PLN/GLN psicologicamente motivadas.

### 3.4 Avaliação e disseminação

Com relação à questão do reconhecimento de personalidade, a avaliação dos modelos produzidos nas atividades (2a..2c) será realizada com uso de medidas comumente adotadas para avaliação de modelos de aprendizagem, tais como precisão, cobertura e medida-F. Além disso, poderá também ser considerada a comparação com a abordagem independente de língua proposto em (Celli, 2012) caso se demonstre viável para o Português<sup>1</sup>.

Com relação à geração de texto, a avaliação dos modelos produzidos (atividades 3a..3d) será realizada de forma intrínseca com base em métricas consagradas da área, como Dice (Dice, 1945) para o caso de conteúdo semântico como em (Paraboni e Deemter, 2014; Silva e Paraboni, 2015), ou BLEU (Papineni et al., 2002) para formas superficiais.

---

<sup>1</sup>O sistema proposto em (Celli, 2012) foi testado apenas nos idiomas inglês e italiano.

Todos os resultados da pesquisa e os produtos desenvolvidos ao longo do projeto serão divulgados em eventos e publicações científicas da área de PLN/GLN e afins.

## 3.5 Recursos

A instituição de destino oferecerá o espaço físico e recursos computacionais básicos (inclusive provenientes de apoio anterior desta e outras agências de fomento) que se façam necessários para o início deste estudo. Os itens solicitados são assim destinados à modernização dos recursos computacionais existentes e participação em reuniões e eventos científicos de divulgação da pesquisa realizada ao longo dos dois anos previstos para a sua execução.

O custo total em material permanente para desenvolvimento do projeto é estimado em R\$ 19.500,00, conforme relacionado a seguir. Além destes recursos, é solicitada também uma bolsa de treinamento técnico de nível superior (TT) para auxílio a diversas atividades ao longo do projeto. As atividades previstas para este bolsista encontram-se anexadas à documentação submetida a esta agência.

1. Três computadores do tipo *desktop* para desenvolvimento do projeto (R\$ 4.000,00 cada, totalizando R\$ 12.000,00).
2. Um computador do tipo *laptop* para gerenciamento e desenvolvimento do projeto (R\$ 4.500,00).
3. Dois dispositivos *nobreak* de 1800VA para segurança dos equipamentos (R\$ 1.500,00 cada, totalizando R\$ 3.000,00).

O projeto será conduzido pelo seu proponente com o apoio de cinco estudantes de mestrado já selecionados, que dedicarão seus projetos individuais às diversas questões de pesquisa discutidas nas seções anteriores.

Além desta equipe, o projeto prevê a possibilidade de colaboração com outros pesquisadores em apoio a diversas tarefas aqui elencadas. No âmbito da própria USP-EACH, destacamos a possibilidade de colaboração com a Dra. Ariane Machado Lima para questões relativas ao uso de técnicas de aprendizagem de máquina, e do Dr. Luciano Antônio Digiampietri para questões ligadas à análise de redes sociais. Destacamos também uma possível colaboração, no âmbito da USP São Carlos, com a Dra. Sandra Maria Aluísio para questões relativas ao uso de conhecimento lexical e psicolinguístico.

### 3.5.1 Outros apoios

De forma complementar à presente solicitação, uma proposta de maior escopo (de três anos de duração) foi submetida à agência CNPq (edital Universal 2016), e encontra-se pendente de resposta.

## Referências

- ALLPORT, F. H.; ALLPORT, G. W. Personality traits: Their classification and measurement. *Journal of Abnormal And Social Psychology*, v. 16, p. 6–40, 1921.
- ANDRADE, J. M. de. *Evidências de validade do inventário dos cinco grandes fatores de personalidade para o Brasil*. Tese (Doutorado) — Universidade de Brasília, 2008.
- ARGAMON, S.; DHAWLE, S.; KOPPEL, M.; PENNEBAKER, J. W. Lexical predictors of personality type. In: *The joint annual meeting of the interface and the classification society of North America*. 2005.
- BATES, J.; NEVILLE, J.; TYLER, J. Using latent communication styles to predict individual characteristics. In: *Proceedings of the 3rd Workshop on Social Media Analytics*. 2012.
- BICK, E. *The Parsing System Palavras - Automatic Grammatical Analysis of Portuguese in a Constraint Grammar Framework*. Tese (Doutorado) — Aarhus University, 2000.
- CALAIS, L. L.; LIMA-GREGIO, A. M.; ARANTES, P.; GIL, D.; BORGES, A. C. L. de C. Um julgamento de concreitude de palavras. *Jornal da Sociedade Brasileira de Fonoaudiologia*, v. 24, p. 262–268, 00 2012.
- CELLI, F. *Adaptive Personality Recognition from Text*. Tese (Doutorado) — University of Trento, 2012.
- DAN-GLAUSER, E. S.; SCHERER, K. R. The Geneva affective picture database (GAPED): a new 730-picture database focusing on valence and normative significance. *Behavior Research Methods*, v. 43, n. 2, p. 468–477, 2011.
- DICE, L. R. Measures of the amount of ecologic association between species. *Ecology*, v. 26, n. 3, p. 297–302, 1945.
- DURAN, M. S.; NUNES, M. G. V. A importância dos falsos homógrafos para a correção automática de erros ortográficos em Português. In: *STIL-2015 IV Jornada de Descrição do Português - 2015*.
- FERREIRA, T. C.; PARABONI, I. Referring expression generation: taking speakers’ preferences into account. *Lecture Notes in Artificial Intelligence*, Springer, v. 8655, p. 539–546, 2014.
- FILHO, P. P. B.; ALUÍSIO, S. M.; PARDO, T. An evaluation of the Brazilian Portuguese LIWC dictionary for sentiment analysis. In: *9th Brazilian Symposium in Information and Human Language Technology - STIL*. 2013. p. 215–219.
- FONSECA, E. R.; ROSA, J. L. G. Mac-Morpho revisited: Towards robust part-of-speech tagging. In: *9th Brazilian Symposium in Information and Human Language Technology*. 2013. p. 98–107.
- GOLDBERG, L. R. An alternative description of personality: The Big-Five factor structure. *Journal of Personality and Social Psychology*, v. 59, p. 1216–1229, 1990.
- IACOBELLI, F.; GILL, A. J.; NOWSON, S.; OBERLANDER, J. Large scale personality classification of bloggers. In: D’MELLO, S. K.; GRAESSER, A. C.; SCHULLER, B.; MARTIN, J.-C. (Ed.). *ACII (2)*. : Springer, 2011. (Lecture Notes in Computer Science, v. 6975), p. 568–577. ISBN 978-3-642-24570-1.
- JANCZURA, G. A.; CASTILHO, G. M. de; ROCHA, N. O.; ERVEN, T. de Jesus Cordeiro van; HUANG, T. P. Normas de concreitude para 909 palavras da língua portuguesa. *Psicologia: Teoria e Pesquisa*, Scielo, v. 23, p. 195–204, 06 2007. ISSN 0102-3772.
- JOHN, O. P.; DONAHUE, E.; KENTLE, R. *The Big Five Inventory - Versions 4a and 54*. Berkeley, CA, USA, 1991.
- JOHN, O. P.; NAUMANN, L. P.; SOTO, C. J. Paradigm shift to the integrative Big-Five trait taxonomy: History, measurement, and conceptual issues. In: \_\_\_\_\_. *Handbook of personality: Theory and research*. New York, NY: Guilford Press, 2008. p. 114–158.
- KOSINSKI, M.; MATZ, S.; GOSLING, S.; POPOV, V.; STILLWELL, D. Facebook as a social science research tool: Opportunities, challenges, ethical considerations and practical guidelines. *American Psychologist*, v. 70, n. 6, p. 543–556, 2015.
- KRAHMER, E.; DEEMTER, K. van. Computational generation of referring expressions: A survey. *Computational Linguistics*, v. 38, n. 1, p. 173–218, 2012.
- LUCENA, D. J. de; PARABONI, I.; PEREIRA, D. B. From semantic properties to surface text: The generation of domain object descriptions. *Inteligencia Artificial. Revista Iberoamericana de Inteligencia Artificial*, Asociacion Espanhola para la Inteligencia Artificial, v. 14, n. 45, p. 48–58, 2010.
- MAIRESSE, F.; WALKER, M.; MEHL, M.; MOORE, R. Using linguistic cues for the automatic recognition of personality in conversation and text. *Journal of Artificial Intelligence Research (JAIR)*, v. 30, p. 457–500, 2007.

- MAIRESSE, F.; WALKER, M. A. Towards personality-based user adaptation: psychologically informed stylistic language generation. *User Model. User-Adapt. Interaction*, v. 20, n. 3, p. 227–278, 2010.
- MAIRESSE, F.; WALKER, M. A. Controlling user perceptions of linguistic style: Trainable generation of personality traits. *Computational Linguistics*, v. 37, n. 3, p. 455–488, 2011.
- MARQUARDT, J.; FARNADI, G.; VASUDEVAN, G.; MOENS, M.-F.; DAVALOS, S.; TEREDESAL, A.; COCK, M. de. Age and gender identification in social media. In: *Proceedings of CLEF 2014 Evaluation Labs, CLEF 2014 Evaluation Labs*. Sheffield: , 2014. p. 1129–1136.
- MUNIZ, M. C. M. *A construção de recursos linguístico-computacionais para o Português do Brasil: o projeto de Uniter-PB*. Dissertação (Mestrado) — ICMC / USP São Carlos, 2004.
- NOVAIS, E. M. de; PARABONI, I. Portuguese text generation using factored language models. *Journal of the Brazilian Computer Society*, Springer-Verlag, v. 19, n. 2, p. 135–146, 2012.
- NOWSON, S.; OBERLANDER, J. Identifying more bloggers: Towards large scale personality classification of personal weblogs. In: *Proceedings of the International Conference on Weblogs and Social Media*. 2007.
- NUNES, M. A. S. N.; TELES, F. R.; SOUZA, J. Inferindo personalidade via tweets. *GEINTEC - Gestão, Inovação e Tecnologias*, v. 3, n. 3, p. 45–57, 2013.
- OBERLANDER, J.; NOWSON, S. Whose thumb is it anyway? classifying author personality from weblog text. In: *COLING/ACL 2006 Poster Sessions*. Sydney, Australia: , 2006. p. 627–634.
- OLIVEIRA, N. R. de; JANCZURA, G. A.; CASTILHO, G. M. de. Normas de alerta e valência para 908 palavras da Língua Portuguesa. *Psicologia: Teoria e Pesquisa*, Scielo, v. 29, p. 185–200, 06 2013. ISSN 0102-3772.
- PAIVA, D. S.; EVANS, R. A framework for stylistically controlled generation. In: BELZ, A.; EVANS, R.; PIWEK, P. (Ed.). *Natural Language Generation*. : Springer Berlin Heidelberg, 2004, (Lecture Notes in Computer Science, v. 3123). p. 120–129.
- PAPINENI, S.; ROUKOS, T.; WARD, W.; ZHU, W. Bleu: a method for automatic evaluation of machine translation. In: *Proceedings of ACL-2002*. 2002. p. 311–318.
- PARABONI, I.; DEEMTER, K. van. Reference and the facilitation of search in spatial domains. *Language, Cognition and Neuroscience*, v. 29, n. 8, p. 1002–1017, 2014.
- PARABONI, I.; DEEMTER, K. van; MASTHOFF, J. Generating referring expressions: Making referents easy to identify. *Computational Linguistics*, MIT Press, v. 33, n. 2, p. 229–254, 2007.
- PARABONI, I.; GALINDO, M.; IACOVELLI, D. Stars2: a corpus of object descriptions in a visual domain. *Language Resources and Evaluation*, Springer, 2016.
- QIU, L.; LIN, H.; RAMSAY, J.; YANG, F. You are what you tweet: Personality expression and perception on Twitter. *Journal of Research in Personality*, v. 46, n. 6, p. 710–718, 2012.
- REITER, E.; DALE, R. *Building natural language generation systems*. New York, NY, USA: Cambridge University Press, 2000. ISBN 0-521-62036-8.
- REITER, E.; WILLIAMS, S. Generating texts in different styles. In: ARGAMON, S.; BURNS, K.; DUBNOV, S. (Ed.). *The Structure of Style*. : Springer Berlin Heidelberg, 2010. p. 59–75. ISBN 978-3-642-12336-8.
- RIGHI, G.; PEISSIG, J. J.; TARR, M. J. Recognizing disguised faces. *Visual Cognition*, v. 20, n. 2, p. 143–169, 2012.
- SCARTON, C. E. Avaliação da inteligibilidade de textos para o público infantil: adaptação das métricas do Coh-Metrix para o Português. In: *Proceedings of STIL 2009*. 2009.
- SCHWARTZ, H. A.; EICHSTAEDT, J. C.; KERN, M. L.; DZIURZYNSKI, L.; RAMONES, S. M.; AGRAWAL, M.; SHAH, A.; KOSINSKI, M.; STILLWELL, D.; SELIGMAN, M.; UNGAR, L. H. Personality, gender, and age in the language of social media: the open-vocabulary approach. *PLoS One*, v. 8, n. 9, p. e73791, 2013.
- SILVA, D. dos S.; PARABONI, I. Generating spatial referring expressions in interactive 3D worlds. *Spatial Cognition & Computation*, v. 15, n. 03, p. 186–225, 2015.
- STENT, A.; MOLINA, M. Evaluating automatic extraction of rules for sentence plan construction. In: *SIGDIAL '09 Proceedings*. 2009. p. 290–297.
- TAUSCZIK, Y. R.; PENNEBAKER, J. W. The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language and Social Psychology*, v. 29, n. 1, p. 24–54, 2010.
- WALKER, M.; STENT, A.; MAIRESSE, F.; PRASAD, R. Individual and domain adaptation in sentence planning for dialogue. *Journal of Artificial Intelligence Research (JAIR)*, v. 30, p. 413–456, 2007.
- YARKONI, T. Personality in 100,000 words: A large-scale analysis of personality and word use among bloggers. *Journal of Research in Personality*, v. 44, n. 3, p. 363–373, 2010. ISSN 0092-6566.