

PGCOMP - Programa de Pós-Graduação em Ciência da Computação
Universidade Federal da Bahia (UFBA)
Av. Milton Santos, s/n - Ondina
Salvador, BA, Brasil, 40170-110

<https://pgcomp.ufba.br>
pgcomp@ufba.br

Resumo: O aumento no número de aplicações que demandam acessibilidade, recuperação de informação e interação humano-computador vem culminando com uma crescente necessidade de geração automatizada da descrição de uma imagem. Essa descrição automatizada requer uma identificação do cenário, dos personagens e dos objetos presentes e de como esses elementos se relacionam entre si. A partir destes elementos torna-se possível gerar uma sentença em linguagem natural descrevendo o conteúdo da imagem. O desenvolvimento de métodos capazes de gerar de uma maneira automática uma sentença que descreve uma imagem permeia uma área de pesquisa denominada *Image Captioning*. A maioria das pesquisas e *datasets* da área de *Image Captioning* se concentram na língua inglesa, desenvolvendo modelos e construindo recursos eficientes no estado da arte. Línguas com poucos recursos para desenvolvimento, tais como o Português, demandam maior pesquisa para alcançarem uma sentença descritiva e compreensível. Porém, somente a aglomeração de vários objetos contidos na imagem não gera uma sentença descritiva de uma cena. Diante deste contexto, este trabalho propõe a análise e incorporação de recursos linguísticos que possam guiar o modelo de linguagem na geração de uma descrição que seja mais informativa da imagem em Português. Experimentos foram realizados com a tradução de *datasets* para a geração da descrição em Português. Os resultados obtidos dão indícios de que existe aprendizado morfológico no treinamento de um modelo de *Image Captioning* e que a incorporação de classes gramaticais durante o treinamento pode contribuir para a geração de sentenças com maior comprimento e mais informativas.

Abstract: The increase in the number of applications that require accessibility, information retrieval and human-computer interaction has culminated in a growing need for automated generation of the description of an image. This automated description requires an identification of the scenario, characters and objects present and how these elements relate to each other. From these elements it becomes possible to generate a sentence in natural language describing the content of the image. The development of methods capable of automatically generating a sentence that describes an image permeates a research area called Image Captioning. Most research and datasets in the Image Captioning area focus on the English language, developing models and building efficient state-of-the-art resources. Languages with few resources for development, such as Portuguese, require more research to achieve a descriptive and understandable sentence. However, only the agglomeration of several objects contained in the image does not generate a descriptive sentence of a scene. In this context, this work proposes the analysis and incorporation of linguistic resources that can guide the language model in generating a description that is more informative of the image in Portuguese. Experiments were performed with the translation of datasets for the generation of the description in Portuguese. The results give evidence that there is morphological learning in the training of an Image Captioning model and that the incorporation of grammatical classes during training can contribute to the generation of sentences with greater length and more informative.

Palavras-chave: Descrição de Imagens; Redes Neurais; Visão Computacional; Processamento de Linguagem Natural. Image Captioning; Neural Networks; Computer Vision; Natural Language Processing.

Ampliando modelos de Image Captioning em Português através das Informações Linguísticas

João Medrado Gondim

Dissertação de Mestrado

Universidade Federal da Bahia

Programa de Pós-Graduação em
Ciência da Computação

Agosto | 2023

MSC | 163 | 2023

Ampliando modelos de Image Captioning em Português através das Informações Linguísticas

João Medrado Gondim

UFBA





Universidade Federal da Bahia
Instituto de Computação

Programa de Pós-Graduação em Ciência da Computação

**AMPLIANDO MODELOS DE IMAGE
CAPTIONING EM PORTUGUÊS ATRAVÉS
DAS INFORMAÇÕES LINGUÍSTICAS**

João Medrado Gondim

DISSERTAÇÃO DE MESTRADO

Salvador
23 de Agosto de 2023

JOÃO MEDRADO GONDIM

**AMPLIANDO MODELOS DE IMAGE CAPTIONING EM
PORTUGUÊS ATRAVÉS DAS INFORMAÇÕES LINGUÍSTICAS**

Esta Dissertação de Mestrado foi apresentada ao Programa de Pós-Graduação em Ciência da Computação da Universidade Federal da Bahia, como requisito parcial para obtenção do grau de Mestre em Ciência da Computação.

Orientadora: Daniela Barreiro Claro
Co-orientador: Marlo Vieira dos Santos e Souza

Salvador
23 de Agosto de 2023

Sistema de Bibliotecas - UFBA

G637 Gondim, João Medrado.

Ampliando modelos de Image Captioning em Português através das Informações Linguísticas / João Medrado Gondim – Salvador, 2023.

76p.: il.

Orientadora: Prof. Dr. Daniela Barreiro Claro.

Co-orientador: Prof. Dr. Marlo Vieira dos Santos e Souza.

Dissertação (Mestrado) – Universidade Federal da Bahia, Instituto de Computação, 2023.

1. Imagens. 2. Redes Neurais. 3. Visão Computacional. 4. Linguagem Natural. I. Claro, Daniela Barreiro. II. Souza, Marlo Vieira dos Santos e. III. Universidade Federal da Bahia. Instituto de Computação. IV. Título.

CDU – 004.5



“Ampliando modelos de Image Captioning com Informações Linguísticas”

JOAO MEDRADO GONDIM

Dissertação apresentada ao Colegiado do Programa de Pós-Graduação em Ciência da Computação na Universidade Federal da Bahia, como requisito parcial para obtenção do Título de Mestre em Ciência da Computação.

Banca Examinadora

Profa. Dra. Daniela Barreiro Claro (Orientadora PGCOMP)

Profa. Dra. Tatiane Nogueira Rios (UFBA-PGCOMP)

Profa. Dra. Sandra Eliza Fontes de Avila (Unicamp)

Dedico este trabalho a quem faz pesquisa no Brasil.

AGRADECIMENTOS

Quero agradecer minha mãe, minha irmã e meu pai pelo apoio na busca pelo tão sonhado mestrado. Nos momentos de maiores pressões e exigências impostas pela pós-graduação foi o apoio e orgulho deles que me mantiveram forte, com meu objetivo final em mente. Às pessoas que estiveram presentes em minha rede de apoio ao longo desses anos, mantendo-me animado em momentos difíceis, meu “muito obrigado”.

Agradeço também a todas e todos do grupo FORMAS, que me ajudaram a entender que pesquisa e ciência são muito melhores quando feitas em grupo. Agradeço minha orientadora, professora Daniela Claro, pela paciência em me guiar e pelos puxões de orelha necessários durante essa jornada.

Por fim, agradeço à Universidade Federal da Bahia pela oportunidade de fazer pós-graduação em uma universidade pública, gratuita e de qualidade. Meu futuro profissional será de eterna gratidão e reconhecimento.

Faz teu melhor e confia.

—BRAZA

RESUMO

O aumento no número de aplicações que demandam acessibilidade, recuperação de informação e interação humano-computador vem culminando com uma crescente necessidade de geração automatizada da descrição de uma imagem. Essa descrição automatizada requer uma identificação do cenário, dos personagens e dos objetos presentes e de como esses elementos se relacionam entre si. A partir destes elementos torna-se possível gerar uma sentença em linguagem natural descrevendo o conteúdo da imagem. O desenvolvimento de métodos capazes de gerar de uma maneira automática uma sentença que descreve uma imagem permeia uma área de pesquisa denominada *Image Captioning*. A maioria das pesquisas e *datasets* da área de *Image Captioning* se concentram na língua inglesa, desenvolvendo modelos e construindo recursos eficientes no estado da arte. Línguas com poucos recursos para desenvolvimento, tais como o Português, demandam maior pesquisa para alcançarem uma sentença descritiva e compreensível. Porém, somente a aglomeração de vários objetos contidos na imagem não gera uma sentença descritiva de uma cena. Diante deste contexto, este trabalho propõe a análise e incorporação de recursos linguísticos que possam guiar o modelo de linguagem na geração de uma descrição que seja mais informativa da imagem em Português. Experimentos foram realizados com a tradução de *datasets* para a geração da descrição em Português. Os resultados obtidos dão indícios de que existe aprendizado morfológico no treinamento de um modelo de *Image Captioning* e que a incorporação de classes gramaticais durante o treinamento pode contribuir para a geração de sentenças com maior comprimento e mais informativas.

Palavras-chave: Descrição de Imagens; Redes Neurais; Visão Computacional; Processamento de Linguagem Natural.

ABSTRACT

The increase in the number of applications that require accessibility, information retrieval and human-computer interaction has culminated in a growing need for automated generation of the description of an image. This automated description requires an identification of the scenario, characters and objects present and how these elements relate to each other. From these elements it becomes possible to generate a sentence in natural language describing the content of the image. The development of methods capable of automatically generating a sentence that describes an image permeates a research area called Image Captioning. Most research and datasets in the Image Captioning area focus on the English language, developing models and building efficient state-of-the-art resources. Languages with few resources for development, such as Portuguese, require more research to achieve a descriptive and understandable sentence. However, only the agglomeration of several objects contained in the image does not generate a descriptive sentence of a scene. In this context, this work proposes the analysis and incorporation of linguistic resources that can guide the language model in generating a description that is more informative of the image in Portuguese. Experiments were performed with the translation of datasets for the generation of the description in Portuguese. The results give evidence that there is morphological learning in the training of an Image Captioning model and that the incorporation of grammatical classes during training can contribute to the generation of sentences with greater length and more informative.

Keywords: Image Captioning; Neural Networks; Computer Vision; Natural Language Processing.

SUMÁRIO

Capítulo 1—Introdução	1
Capítulo 2—Introdução	3
2.1 Motivação	4
2.2 Definição do Problema	5
2.3 Objetivos	5
2.4 Questões de pesquisa	6
2.5 Organização do documento	7
Capítulo 3—Fundamentação Teórica	9
3.1 Redes Neurais	9
3.1.1 Rede Neural Convolutacional	11
3.1.2 Redes Neurais Recorrentes	12
3.1.3 Modelos de Linguagem	13
3.1.4 Arquiteturas Encoder-decoder	14
3.1.5 Mecanismo de Atenção	15
3.1.6 Transformers	17
3.1.7 Evolução dos Métodos	19
3.2 Processamento de Linguagem Natural	21
3.2.1 Análise Morfológica	22
3.2.2 Análise Sintática	24
3.2.2.1 Classes Gramaticais	25
3.3 Métricas de Avaliação	26
3.3.1 BiLingual Evaluation Understudy - BLEU	27
3.3.2 Metric for Evaluation of Translation with Explicit Odering - METEOR	29
3.3.3 Consensus-based Image Description Evaluation - CIDEr	30
3.3.4 Semantic Propositional Image Caption Evaluation - SPICE	31
3.4 Resumo do Capítulo	32
Capítulo 4—Trabalhos Relacionados	33
4.1 Image captioning	33
4.1.1 Abordagens Iniciais	33
4.1.2 Image captioning com Redes Neurais	34

4.1.2.1	CNN → RNN	34
4.1.2.2	CNN → <i>Attention</i> → RNN	35
4.1.2.3	CNN → Transformer	37
4.2	Conjuntos de Dados	40
4.3	Resumo do Capítulo	44
Capítulo 5—Linguistic Image Captioning		45
5.1	Metodologia	45
5.1.1	Conjuntos de Dados	45
5.1.2	Extração de Classes gramaticais	47
5.1.3	Avaliação	48
5.2	Arquitetura Proposta	49
5.3	Experimentos	51
5.3.1	Experimento 1	51
5.3.2	Resultados 1	52
5.3.3	Experimento 2	53
5.3.4	Resultados 2	54
5.3.5	Experimento 3	55
5.3.6	Resultados 3	58
5.3.7	Experimento 4	62
5.3.8	Resultados 4	63
5.4	Resumo do Capítulo	65
Capítulo 6—Conclusão e Trabalhos Futuros		67
6.1	Publicações	68
6.2	Limitações	69
6.3	Trabalhos Futuros	69
Referências Bibliográficas		71

LISTA DE FIGURAS

3.1	Modelo de um neurônio. Fonte: (RUSSELL; NORVIG, 2020)	10
3.2	Rede perceptron com duas entradas e duas saídas (a) e rede neural com duas entradas, uma camada oculta e duas unidades de saída (b). Fonte: adaptado de (RUSSELL; NORVIG, 2020).	10
3.3	Exemplo de funcionamento do <i>kernel</i> em uma entrada em formato <i>grid</i> . Fonte: adaptado de (ZHANG et al., 2021a).	11
3.4	Esquema de uma Rede Neural Convolutacional. Fonte: (SULTANA; SUFIAN; DUTTA, 2019).	12
3.5	Esquema de uma Rede Neural Recorrente. Fonte: (JURAFSKY; MARTIN, 2022)	12
3.6	Modelo de linguagem utilizando RNN. Fonte: adaptado de (ZHANG et al., 2021a)	13
3.7	Arquitetura <i>Encoder-decoder</i> . Fonte: (JURAFSKY; MARTIN, 2022) . . .	14
3.8	Exemplo de tradução da sentença “They are watching.” para “Eles estão assistindo.”. Fonte: adaptado de (ZHANG et al., 2021a)	15
3.9	Arquitetura <i>Encoder-decoder</i> . Fonte: (JURAFSKY; MARTIN, 2022) . . .	15
3.10	Exemplo de arquitetura <i>Encoder-decoder</i> com uso de atenção. Fonte: adaptado de (JURAFSKY; MARTIN, 2022)	16
3.11	Arquitetura <i>Encoder-decoder</i> de um Transformer. Fonte: adaptado de (VASWANI et al., 2017)	18
3.12	Arquitetura <i>Encoder-decoder</i> aplicada a <i>Image captioning</i> com uma CNN extraindo representações de uma imagem e gerando o vetor de contexto. Fonte: adaptado de (STEFANINI et al., 2021)	19
3.13	Arquitetura <i>Encoder-decoder</i> aplicada a <i>Image captioning</i> com uma CNN extraindo representações de uma imagem e utilizando atenção para gerar um vetor de contexto apropriado para cada palavra. Fonte: adaptado de (STEFANINI et al., 2021)	20
3.14	Arquitetura <i>Encoder-decoder</i> aplicada a <i>Image captioning</i> com uma CNN extraindo representações de uma imagem e utilizando <i>Transformers</i> para criar dependências globais e gerar uma descrição para a entrada. Fonte: adaptado de (STEFANINI et al., 2021) e (VASWANI et al., 2017)	21
3.15	Diagrama com as análises que podem ser feitas de uma sentença. Fonte: (SILVA, 2010)	23
3.16	Árvore sintática para a frase “Eu prefiro comer feijão com arroz”. Fonte: preparado pelos autores.	25
3.17	Imagem de exemplo para geração de grafo de cena. Fonte: (ANDERSON et al., 2016).	31

4.1	Exemplo de descrição de imagem com CNN e RNN em (VINYALS et al., 2015)	34
4.2	Arquitetura do modelo em (HE et al., 2017). Fonte: (HE et al., 2017) . .	35
4.3	Arquitetura do modelo em (ZHANG et al., 2021b). Fonte: (ZHANG et al., 2021b)	36
4.4	Exemplo de geração de descrição em (XU et al., 2016). Fonte: (XU et al., 2016)	37
4.5	Exemplo de geração de descrição em (LU et al., 2017) onde a informação visual é utilizada ou não a depender a palavra predita.	37
4.6	Arquitetura do modelo em (HUANG et al., 2019). Fonte: (HUANG et al., 2019)	38
4.7	Arquitetura do bloco de atenção unificado utilizando <i>pooling linear</i> . Fonte: (PAN et al., 2020)	39
4.8	Arquitetura do modelo proposto em (WANG et al., 2022). Fonte: (WANG et al., 2022)	40
4.9	Exemplo de imagem do Flickr8k. Fonte: (HODOSH; YOUNG; HOCKENMAIER, 2013)	41
4.10	Exemplo de imagem do <i>Microsoft COCO Dataset</i> . Fonte: (LIN et al., 2015)	42
4.11	Exemplo de imagem do # PraCegoVer. Fonte: (SANTOS; COLOMBINI; AVILA, 2022)	43
5.1	Exemplo de imagem do Flickr8k. Fonte: (HODOSH; YOUNG; HOCKENMAIER, 2013)	46
5.2	Exemplo de imagem do <i>Flickr30k</i> . Fonte: (YOUNG et al., 2014)	46
5.3	Arquitetura do modelo de descrição linguística de imagens do trabalho Fonte: adaptado de (VASWANI et al., 2017).	50
5.4	Arquitetura simplificada do modelo de descrição de imagens utilizado em (GONDIM.; CLARO.; SOUZA., 2022) Fonte: preparado pelos autores. .	51
5.5	Comparações entre métricas de descrições feitas com as mesmas imagens e linguagens distintas.	53
5.6	Contagens de marcações feitas nos formulários. Fonte: preparado pelos autores.	54
5.7	Exemplo de predição de sentença e do “foco” dado pela <i>head</i> 1.	56
5.8	Médias de cada <i>head</i> de atenção quando a palavra “andando” foi predita pelo Modelo-PT.	57
5.9	Médias de cada <i>head</i> de atenção quando a palavra “ <i>walking</i> ” foi predita pelo Modelo-EN.	57
5.10	Médias de cada <i>head</i> de atenção quando a palavra “andar” foi predita pelo Modelo-PT.	57
5.11	Diferenças entre as <i>heads</i> do Modelo-PT ao predizer as palavras “andando” e “andar”.	58
5.12	Comparações entre diferenças de pesos de atenção retornados pelo Modelo-PT. Acima, o mesmo verbo com sufixos diferentes. Embaixo, dois verbos diferentes, porém com o mesmo sufixo.	59

5.13	Comparações entre diferenças de pesos de atenção retornados pelo Modelo-PT. Acima, o mesmo verbo com sufixos diferentes. Embaixo, dois verbos diferentes, porém com o mesmo sufixo.	59
5.14	Diferenças entre as <i>heads</i> do Modelo-PT ao predizer as palavras “branca” e “branco”.	60
5.15	Diferenças entre as <i>heads</i> do Modelo-PT ao predizer as palavras “vermelha” e “vermelho”.	60
5.16	Comparações entre diferenças de pesos de atenção retornados pelo Modelo-PT. Acima, o mesmo adjetivo com sufixos diferentes. Embaixo, dois adjetivos diferentes também com sufixos diferentes.	61
5.17	Diferenças entre as <i>heads</i> do Modelo-PT ao predizer as palavras “homem” e “mulher”.	61
5.18	Comparações entre diferenças de pesos de atenção retornados pelo Modelo-PT.	62

LISTA DE TABELAS

3.2	Tabela com os rótulos e significados de cada POS de acordo com (RADE-MAKER et al., 2017).	26
4.1	Tabela com os dados dos <i>datasets</i> COCO, Flickr8k, <i>#PraCegoVer-63K</i> e <i>#PraCegoVer-173K</i> . Fonte: adaptado de (SANTOS; COLOMBINI; AVILA, 2022).	43
5.1	Tabela com informações dos conjuntos de dados utilizados nos experimentos.	47
5.2	Métricas obtidas com o modelo treinado com o conjunto de dados em Português.	52
5.3	Métricas obtidas do modelo treinado com o conjunto de dados em Inglês.	52
5.4	Pontuações de concordâncias entre anotadores para cada erro.	55
5.5	Métricas obtidas para cada modelo treinado e avaliado com seus respectivos conjuntos de dados.	58
5.6	Métricas obtidas para cada modelo treinado e avaliado com seus respectivos conjuntos de dados utilizando o modelo proposto (Modelo-LIC).	63

LISTA DE SIGLAS

PLN	Processamento de Linguagem Natural	3
VC	Visão Computacional	3
FFN	Feedforward Neural Networks	10
CNN	Convolutional Neural Network	11
RNN	Recurrent Neural Network	12
NLG	Natural Language Generation	22
NLU	Natural Language Understanding	22
POS	Part-of-speech	25
BLEU	Bilingual Evaluation Understudy	26
METEOR	Metric for Evaluation of Translation with Explicit Odering	26
CIDeR	Consensus-based Image Description Evaluation	26
SPICE	Semantic Propositional Image Caption Evaluation	26
LIC	Linguistic Image Captioning	49

Capítulo

1

Este capítulo introduz a área de Image captioning. O início contextualiza o o tópico e depois apresentam-se as motivações, objetivos e hipóteses.

INTRODUÇÃO

Neste capítulo introduzimos a área de Image captioning. Iniciamos contextualizando o tópico e depois apresentando as motivações, os objetivos e hipóteses.

INTRODUÇÃO

Visualizar uma imagem, compreender seus conteúdos e expressar de maneira natural o que foi entendido é uma habilidade natural do ser humano (BAI; AN, 2018). Automatizar esse processo unicamente com um computador, é uma tarefa que compartilha ferramentas de áreas distintas da Inteligência Artificial: Visão Computacional (VC) (do inglês *Computer Vision*) e Processamento de Linguagem Natural (PLN) (do inglês *Natural Language Processing*). De VC são empregados sistemas capazes de entender e identificar o conteúdo de uma imagem, gerando representações visuais do que se deseja descrever. Já de PLN, utiliza-se modelos de linguagem capazes de gerar sentenças a partir da informação extraída (STEFANINI et al., 2021). À tarefa automatizada de criar uma descrição de imagem é dada o nome *Image captioning*. O desenvolvimento de métodos que geram uma descrição para uma imagem têm aumentado nos últimos anos, com pesquisas buscando como compreender os conteúdos da imagem e gerar mais informação do que apenas uma lista de classes ou objetos contidos nela (SHARIF et al., 2020).

A área de *Image captioning* apresenta diferentes aplicações. A acessibilidade de usuários com visão dificultada melhora muito com a presença de texto alternativo em imagens na web (Web Accessibility Initiative, 2022) e a falta de descrições de imagens pode ser suprida com modelos automatizados que cumpram essa função. Muitas imagens existentes na web possuem textos alternativos atrelados e buscadores contam com esses textos para encontrá-las, quando não há texto alternativo, associar corretamente um texto de busca à uma imagem melhora a pesquisa feita (SHARIF et al., 2020). Ademais, descrever automaticamente uma imagem é uma tarefa explorada por organizadores de galerias de fotos, analisando imagens semelhantes e juntando-as em grupos com uma sentença em comum ou sistemas de interação humano computador, auxiliando robôs a interagirem com o ambiente de acordo com suas percepções visuais (SHARIF et al., 2020).

Image captioning combina métodos de duas áreas diferentes, sendo portanto desafiadora tanto do ponto de vista de VC quanto de PLN (SHARIF et al., 2020). Do ponto de vista de VC, identificar corretamente os conteúdos visuais “relevantes” da cena é importante uma vez que uma figura pode ter muitos objetos, mas humanos podem não dar

importância a todos no momento da descrição, portanto precisa-se saber a quais partes serão dadas maior importância ao analisar uma imagem. Já para a área de PLN, modelos de descrição de imagens podem ter uma tendência de memorizar objetos ou palavras que apareçam próximos durante o treinamento (DOGNIN et al., 2018), diminuindo sua capacidade de generalização¹. Além disso, avaliar um modelo de *Image captioning* é uma tarefa difícil por conta da heterogeneidade entre entrada e saída, sendo algumas das métricas automatizadas empregadas para avaliar saídas e comparar modelos ainda adaptadas da tarefa de tradução automática², como a BLEU (PAPINENI et al., 2002) e a METEOR (BANERJEE; LAVIE, 2005). Ainda que sistemas de avaliação tenham sido criados especificamente para a área, CIDER (VEDANTAM; ZITNICK; PARIKH, 2014) e SPICE (ANDERSON et al., 2016), os autores em (MILTENBURG; ELLIOTT, 2017) observam que métricas baseadas em texto somente não são o suficiente para avaliar desvantagens de modelos, sendo importante avaliar qualitativamente as saídas geradas.

2.1 MOTIVAÇÃO

Em (MIELKE, 2016) é apontado que 69% dos artigos publicados no *Annual Meeting of the Association for Computational Linguistics* (ACL) de 2016 tiveram como língua avaliada o Inglês. A maioria das pesquisas em modelos e conjuntos de dados que existem se concentraram na língua inglesa. Apenas o *#PraCegoVer* (SANTOS; COLOMBINI; AVILA, 2022) pode ser citado como um conjunto de dados preparado para treinamento de modelos de *Image captioning* em português, tendo sido publicado em 2022. Em (BENDER, 2019), a autora aponta que existem 240 famílias de linguagens no mundo, com aproximadamente 7000 linguagens derivadas delas, e ainda assim há relatos de artigos que sequer informam a linguagem utilizada nos experimentos. Segundo a autora, a língua inglesa não pode ser utilizada para representar outras linguagens devido a certas especificidades linguísticas como: poucas variações entre palavras por conta de sua “morfologia pequena” ou ordem fixa de palavras. Em (MILTENBURG; ELLIOTT, 2017) os autores categorizam tipos de erros existentes em um sistema estado da arte de *Image captioning* para fornecer um ponto de vista qualitativo ao avaliar as saídas de um modelo. Os autores concluem que avaliações complementares às feitas com métricas automatizadas de texto podem ser consideradas relevantes para mensurar os avanços de sistemas de descrição de imagens, permitindo análises que vão além da similaridade entre sentenças. Considerando a análise qualitativa de erros como passo importante para investigar e criar arquiteturas de modelos de descrição de imagens, as discrepâncias entre as características de linguagens diferentes podem influenciar nessa análise, ou seja, basear-se em erros de um modelo treinado em uma linguagem pode dificultar a generalização da mesma arquitetura treinada em outros idiomas. Portanto, uma das principais motivações deste trabalho é a falta de estudos na área de *Image captioning* para a língua portuguesa.

Além disso, autora em (BENDER, 2009) e (BENDER, 2013) discute que para que um sistema de geração de linguagem funcione independente do idioma, este deve utilizar

¹Habilidade de se adaptar a dados não vistos.

²Gerar texto em um linguagem a partir de uma entrada de texto em uma linguagem diferente (GATT; KRAHMER, 2018).

conhecimentos linguísticos no seu desenvolvimento. Portanto, a falta de informações gramaticais durante o desenvolvimento de um sistema de *Image captioning* pode acarretar em descrições errôneas estruturalmente. A segunda motivação deste trabalho é utilizar características linguísticas para avaliar e ampliar a capacidade de modelos de *Image captioning*, ainda que utilizando um conjunto de dados traduzido.

2.2 DEFINIÇÃO DO PROBLEMA

O problema investigado neste trabalho é a falta de modelos desenvolvidos para *Image captioning* na língua portuguesa. Focar na língua inglesa ao se pesquisar modelos de PLN pode introduzir vieses indesejáveis no decorrer do processo de estudo, uma vez que a língua inglesa pode não ser representativa o suficiente. Diferenças que podem ser fonéticas, morfológicas, sintáticas e semânticas fazem com que o desenvolvimento de sistemas em apenas uma linguagem sejam prejudiciais quando se pensa em aplicações que utilizam modelos de PLN criados com foco voltado para um idioma apenas (BENDER, 2019). Por exemplo, na frase em inglês “*A yellow car on the street*”, primeiro tem-se um adjetivo (*yellow*) que define a cor do substantivo carro (*car*), mas na tradução para o português, “Um carro amarelo na rua”, o substantivo (carro) aparece primeiro e depois é seguido pelo adjetivo (amarelo), em ordens diferentes. Ainda na frase acima, substituindo o objeto “*car*” por “*motorcyle*” (“*A yellow motorcyle on the street*”) apenas uma palavra é alterada, enquanto que, no português, essa mudança de substantivo impacta em ambas as palavras adjacentes: “Uma motocicleta amarela na rua”. No português, há uma concordância de gênero entre palavras, como visto anteriormente, mas isso não acontece ao mudar o substantivo na frase em inglês. Traduzindo a frase em português (“Uma motocicleta amarela na rua”) para inglês, a palavra “Uma” pode ter um sentido ambíguo, sendo um artigo ou numeral, e podemos ter duas traduções corretas: “*A yellow motorcyle on the street*” ou “*One yellow motorcyle on the street*”. Ao gerar as palavras de uma sentença, um sistema de *Image captioning* em português deve se atentar para o gênero do objeto identificado em seu passo visual (da área de VC), adaptando as palavras ao redor do substantivo, isso não é considerado caso o sistema seja feito em inglês.

2.3 OBJETIVOS

Considerando as motivações apresentadas junto com o problema a ser investigado, o presente trabalho tem como objetivo principal:

Desenvolver uma arquitetura para descrever imagens em Português utilizando informações linguísticas que sejam capazes de melhorar a geração de textos.

Destacam-se como objetivos específicos:

- OE 1: Analisar um conjunto de dados traduzido do inglês para ser utilizado no treinamento de um modelo de *Image captioning*;
- OE 2: Avaliar a incorporação de informações linguísticas para a geração de texto em um modelo de *Image captioning*.

2.4 QUESTÕES DE PESQUISA

Com o intuito de analisar o OE 1, elaborou-se duas questões de pesquisa: QP 1 e QP 2 que são detalhadas a seguir.

- QP 1: *O uso de conjuntos de dados traduzidos permite o desenvolvimento de sistemas de Image captioning em outras linguagens?*

No início da pesquisa, a falta de conjuntos de dados com sentenças em português era um problema para estudos na área. Uma possível solução é a tradução de um conjunto de dados para o português, embora traduzir sentenças possa introduzir erros no início do processo (ROSA et al., 2021). O objetivo dessa pergunta é entender quais as possíveis desvantagens derivadas do uso de um conjunto de dados traduzido, sendo abordada na seção 5.3.2.

- QP 2: *Quais tipos de erros estão presentes nas saídas de um sistema de Image captioning?*

Essa pergunta tem como objetivo entender e agrupar as falhas cometidas quando se tenta descrever uma imagem. Categorizar esses erros tem como intuito conseguir trabalhar cada um deles individualmente ou em grupos. O resultados na seção 5.3.4 abordam esta pergunta.

Para avaliar o OE 2, foram elaboradas as questões de pesquisa QP 3 e QP 4.

- QP 3: *Modelos de Image Captioning são capazes de aprender relações morfológicas entre imagem e palavra?*

Antes de considerar a incorporação de informações linguísticas, esta pergunta pretende avaliar se algum tipo de compreensão acerca dessas informações é aprendido no processo de treinamento de uma rede neural criada para a tarefa de descrição de imagens. Essa pergunta é endereçada na seção 5.3.5.

- QP 4: *Auxiliar o treinamento de um sistema de Image Captioning com informação das classes gramaticais de cada palavra melhora a sentença gerada?*

Em (BENDER, 2013), lê-se “Morfofossintaxe é a diferença entre uma sentença e uma *bag-of-words*³”, assim sendo, o objetivo dessa pergunta é saber se as relações existentes na frase podem ajudar a guiar o modelo para o local a ser atendido na imagem durante as fases de treinamento e inferência. Essa questão é abordada na seção 5.3.7.

³*Bag-of-words* é uma técnica de extração de características de um texto em que é descrita apenas a ocorrência de palavras no documento analisado, descartando-se informações acerca da estrutura ou ordem das palavras usadas (BROWNLEE, 2017).

2.5 ORGANIZAÇÃO DO DOCUMENTO

Este trabalho está organizado da seguinte forma: no *Capítulo 2* é apresentada a fundamentação teórica necessária para o entendimento da temática de *Image captioning* e das características linguísticas que são utilizadas para auxiliar no desenvolvimento do sistema; em seguida, no *Capítulo 3*, os trabalhos relacionados na área de descrição de imagens são apresentados; o *Capítulo 4* apresenta os conjuntos de dados que foram utilizados, a extração das informações linguísticas empregadas, como o trabalho será avaliado, a arquitetura proposta, experimentos e resultados alcançados; concluimos no *Capítulo 5* com o estado atual da pesquisa, contribuições feitas, publicações e trabalhos futuros.

FUNDAMENTAÇÃO TEÓRICA

Este capítulo apresenta uma visão geral sobre a área de descrição de imagens. Primeiramente são apresentados os métodos utilizados para extração de características da imagem e como são empregados junto a modelos linguagem na geração de sentenças em arquiteturas do tipo *Encoder-decoder* ou *Sequence-to-sequence*. Em seguida descreve-se o mecanismo de atenção, como seu funcionamento auxilia o sistema de *Image captioning* no direcionamento para partes relevantes da imagem a cada palavra gerada e como ele é utilizado para a construção de *Transformers*, arquitetura presente na maioria dos modelos estado da arte para *Image captioning* (STEFANINI et al., 2021). A segunda parte do capítulo explora as informações linguísticas presentes nesta proposta de trabalho. Por fim, as métricas existentes para avaliar o desempenho de modelos de descrição de imagem são apresentadas.

3.1 REDES NEURAIS

Uma rede neural é um conjunto de unidades conectadas por ligações direcionadas. Uma rede neural é definida pelas propriedades dessas unidades, chamadas de “neurônios”, a Figura 3.1 ilustra um modelo de neurônio. Uma unidade i da rede neural possui uma entrada a_0 atrelada a um peso w_{0j} (chamado de viés) e outras entradas a_i atreladas a pesos w_{ij} , esses pesos determinam a força e o sinal de cada ligação entre unidades i e j . Cada unidade j calcula uma soma (ponderada por cada peso w_{ij}) de suas entradas ($in_j = \sum_{i=0}^n w_{i,j}a_i$) e aplica uma função de ativação g para obter a saída da unidade. O uso de redes neurais pode ser de camada única (também chamada de rede *perceptron*) onde todas as unidades de entrada são conectadas diretamente com as unidades de saída, conforme mostrado na Figura 3.2a em que uma rede neural tem duas entradas conectadas a duas saídas, ou com múltiplas camadas, como na Figuras 3.2b em que, além das unidades de entrada e saída, há uma camada “oculta” com unidades entre elas (RUSSELL; NORVIG, 2020).

Ao uso de muitas camadas ocultas em redes neurais modernas é dado o nome de “aprendizado profundo” (do inglês *deep learning*). Por conta de sua alta capacidade de

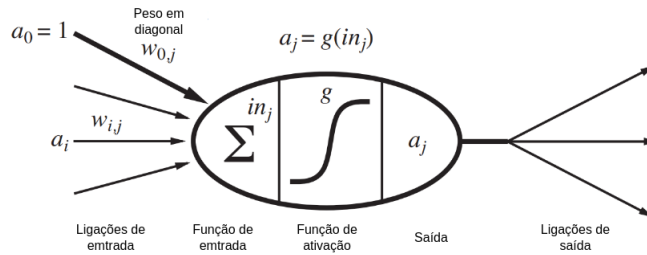


Figura 3.1: Modelo de um neurônio. Fonte: (RUSSELL; NORVIG, 2020)

aprendizado e de paralelização de operações, redes neurais profundas são muito usadas para representação de informação, sendo boas ferramentas para problemas de larga escala, em que o aprendizado automatizado de características da entrada (do inglês *features*) se faz necessário (JURAFSKY; MARTIN, 2022).

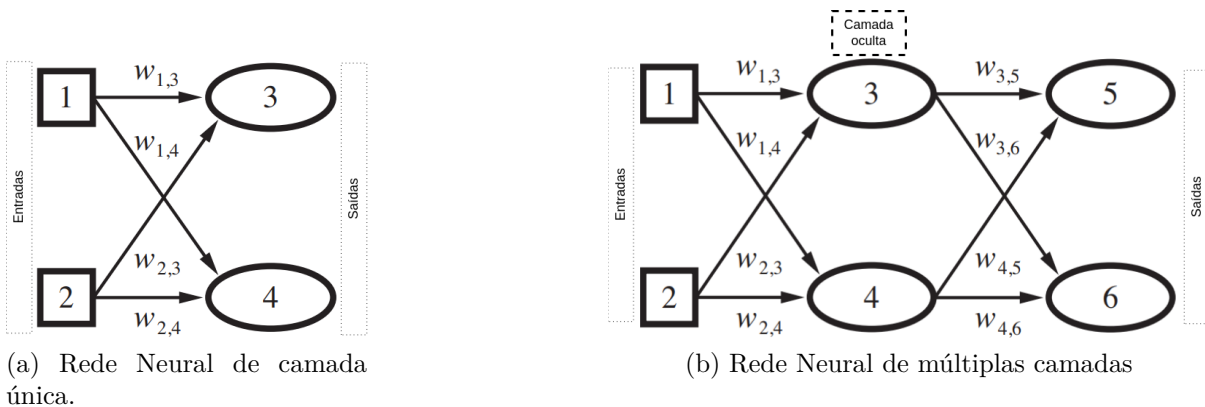


Figura 3.2: Rede perceptron com duas entradas e duas saídas (a) e rede neural com duas entradas, uma camada oculta e duas unidades de saída (b). Fonte: adaptado de (RUSSELL; NORVIG, 2020).

Existem tipos diferentes de Redes Neurais, caracterizadas por suas unidades e como elas se conectam entre si. As redes neurais com conexão para frente (do inglês *Feed-forward Neural Networks (FFN)*) são redes neurais que possuem conexões apenas em uma direção entre suas unidades, formando um grafo acíclico simples (RUSSELL; NORVIG, 2020). Nelas, cada unidade representa a saída de sua entrada atual, não possuindo estado interno. Já uma Rede Neural Recorrente é capaz de armazenar informações de entradas recentes em sua ativação (HOCHREITER; SCHMIDHUBER, 1997). Em uma Rede Neural Convolutiva, o cálculo entre entradas e pesos é feito com o uso da operação de convolução (GOODFELLOW; BENGIO; COURVILLE, 2016).

3.1.1 Rede Neural Convolutacional

As Redes Neurais Convolutacionais (do inglês *Convolutional Neural Network (CNN)*) são “redes neurais que utilizam a operação matemática de convolução no lugar de multiplicação de matrizes em pelo menos uma de suas camadas” (GOODFELLOW; BENGIO; COURVILLE, 2016). As CNNs se tornam excelentes alternativas para dados em formato de grade (“*grid*”)¹ uma vez que são especializadas para este tipo de entrada por três características: interações esparsas, compartilhamento de parâmetros e representações equivalentes. Um exemplo de dado em formato de grade é visto na entrada da Figura 3.3. As interações esparsas de uma CNN diferem das redes neurais convencionais apresentadas anteriormente, em que cada unidade saída de uma camada interage com cada unidade de entrada. Isso é alcançado devido ao uso de um *kernel*² que é geralmente muito menor do que o tamanho da entrada. Também chamadas de conexões esparsas ou pesos esparsos, esse tipo de interação ajuda a capturar características menores e mais significativas em uma imagem por exemplo. Com isso pode-se aumentar a acurácia de uma rede armazenando menos parâmetros, o que economiza espaço de memória (GOODFELLOW; BENGIO; COURVILLE, 2016).

Entrada		Kernel			Saída			
0	1	2	*	=	19	25		
3	4	5					0	1
6	7	8					2	3

Figura 3.3: Exemplo de funcionamento do *kernel* em uma entrada em formato *grid*. Fonte: adaptado de (ZHANG et al., 2021a).

O compartilhamento de parâmetros se refere ao uso compartilhado de alguns parâmetros adaptados do *kernel* em mais de uma localização da entrada, portanto, diferente do que acontece em redes neurais não convolutacionais, há uma redução do armazenamento requerido para esses parâmetros, tornando a convolução mais eficiente que a multiplicação de matrizes (GOODFELLOW; BENGIO; COURVILLE, 2016). O compartilhamento de parâmetros também acarreta em equivariância nas representações, ou seja, as saídas de uma camada convolutacional mudam na mesma proporção que as entradas. Isso é útil para detectar bordas em lugares diferentes de uma imagem (ZHANG et al., 2021a).

A Figura 3.4 mostra um esquema de uma CNN, cada camada convolutacional é seguida de uma camada de *pooling*, que altera as saídas da camada convolutacional por meio de

¹Uma imagem pode ser representada por uma matriz retangular de valores (como uma grade), a esses valores dá-se o nome de *pixel*. Um *pixel* pode indicar, com um número real, o nível de cinza de um elemento da matriz (um ponto na imagem) ou, com uma tripla, os números representando as misturas de cores que resultam na cor final daquele ponto da imagem (HUGHES et al., 2013).

²Um *kernel* é um “vetor multidimensional de parâmetros que são adaptados pelo algoritmo de aprendizado” (GOODFELLOW; BENGIO; COURVILLE, 2016) aplicado à função de convolução da camada da rede.

sínteses estatísticas das saídas próximas. *Pooling* também é chamado de “redução de amostra” (do inglês *downsampling*) pois é utilizado para reduzir a dimensão de saída de uma camada convolucional, permitindo criar CNNs onde a saída da última camada tem uma dimensionalidade muito baixa (AGHDAM; HERAVI, 2017). Os tipos mais comuns de *pooling* são *max-pooling* e *average-pooling* que calculam, respectivamente, o máximo e a média das características das saídas computadas (CARREIRA et al., 2012). Além destes, o *pooling* bilinear (*bilinear pooling*) é capaz de calcular interações de segunda ordem entre as saídas com o uso de produto tensorial (CARREIRA et al., 2012).

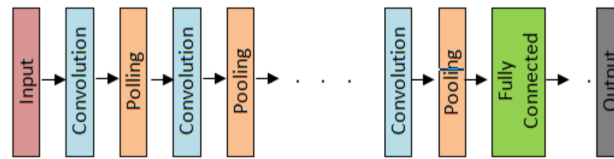


Figura 3.4: Esquema de uma Rede Neural Convolucional. Fonte: (SULTANA; SUFIAN; DUTTA, 2019).

Ao gerar uma sentença para uma imagem, um ser humano precisa primeiro visualizar a imagem a ser descrita. Fazendo uma analogia com a tarefa de *Image captioning*, é necessário “visualizar” a imagem antes de descrevê-la (STEFANINI et al., 2021), portanto, o início do processo consiste em extrair uma representação da imagem. Uma ferramenta adequada para construir uma codificação representativa da imagem é a CNN. Mas, após essa etapa, o conteúdo visual criado precisa ser utilizado para gerar a sentença, para isso emprega-se uma arquitetura de rede neural capaz de gerar uma frase palavra-a-palavra: as Redes Neurais Recorrentes.

3.1.2 Redes Neurais Recorrentes

As Redes Neurais Recorrentes (do inglês *Recurrent Neural Network (RNN)*) são redes neurais especializadas no processamento de dados em forma sequencial (previsão de tempo, curva de temperatura de um paciente, preços de ações ao longo dos dias) (ZHANG et al., 2021a). Diferente de uma CNN ou FFN, a RNN compartilha parâmetros de um estado anterior no estado atual.

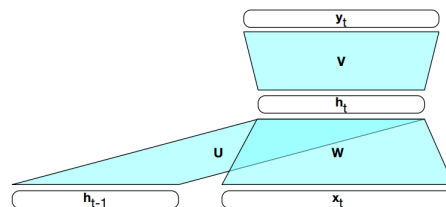


Figura 3.5: Esquema de uma Rede Neural Recorrente. Fonte: (JURAFSKY; MARTIN, 2022)

Um exemplo gráfico de uma RNN é ilustrado na Figura 3.5. Cada estado interno

atual (\mathbf{h}_t) é obtido por uma função g da entrada inserida na rede (\mathbf{x}_t) e do estado anterior (\mathbf{h}_{t-1}), respectivamente multiplicados pelos pesos correspondentes (\mathbf{W} e \mathbf{U}), essa interação é representada na Equação 3.1.2.1 onde cada peso W e U determina o uso da entrada atual x_t e do estado interno anterior h_{t-1} respectivamente. A saída final é calculada pela Equação 3.1.2.2, em que V representa a matriz de pesos entre a camada do estado atual e a camada de saída (JURAFSKY; MARTIN, 2022). Esse compartilhamento de parâmetros é resultante da formulação recorrente da rede, o que torna RNNs capazes de armazenar uma “memória” de entradas anteriores que permanece nos estados internos da rede influenciando as suas saídas futuras (GRAVES, 2012).

$$h_t = g(Uh_{t-1} + Wx_t) \quad (3.1.2.1)$$

$$y_t = f(Vh_t) \quad (3.1.2.2)$$

RNNs podem ser utilizadas para diversas tarefas em Processamento de Linguagem Natural: rotulagem de sequência (*sequence labeling*), classificação de sequência ou modelos de linguagem (JURAFSKY; MARTIN, 2022). Após extrair uma representação de uma imagem, um modelo de arquitetura neural que seja capaz de gerar uma sentença a partir dessa representação extraída é utilizado: um modelo de linguagem.

3.1.3 Modelos de Linguagem

Modelos de linguagem são modelos capazes de prever o próximo elemento de uma sequência de palavras dado um contexto prévio (MIKOLOV et al., 2010). Por exemplo, um modelo de linguagem define a probabilidade de a palavra “mar” vir depois da sequência de palavras “Uma tartaruga no”, ou seja: $P(\text{mar} \mid \text{Uma tartaruga no})$. O uso de RNNs como modelos de linguagem ocorre por sua capacidade preditiva baseada na palavra anteriormente gerada. Como ilustrado na Figura 3.6, por conta do uso do estado interno gerado anteriormente na operação que determina o estado interno atual, uma RNN é capaz de utilizar informações prévias de uma sentença para prever uma palavra futura, essa capacidade de armazenar informações passadas pode preceder até o início da sentença sendo gerada (JURAFSKY; MARTIN, 2022).

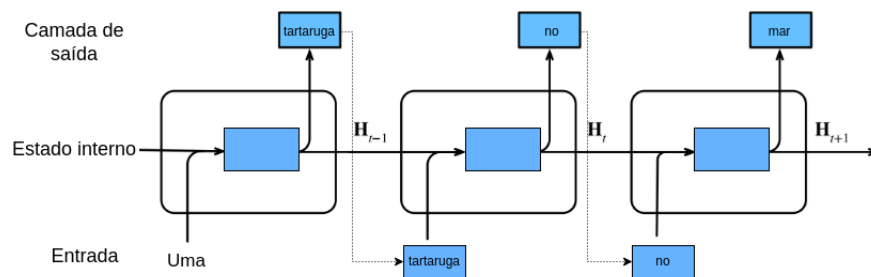


Figura 3.6: Modelo de linguagem utilizando RNN. Fonte: adaptado de (ZHANG et al., 2021a)

No exemplo mostrado na Figura 3.6 a sentença é gerada a partir de uma entrada

única fornecida (“Uma”), mas o uso de RNNs como modelos de linguagem também pode ser empregado em arquiteturas que geram uma sequência a partir de outra sequência de entrada: os modelos *Encoder-decoder*. Os modelos que usam a arquitetura *Encoder-decoder* são utilizados na tarefa de tradução automatizada de sentenças, por aprenderem representações intermediárias entre uma linguagem e outra (CHO et al., 2014b).

3.1.4 Arquiteturas Encoder-decoder

Arquiteturas de redes *Encoder-decoder* (ou *sequence-to-sequence*) são capazes de, partindo de uma sequência de entrada, gerar uma sequência de saída que não necessariamente precisa ter o mesmo comprimento da entrada inicial, mas que é contextualmente apropriada (JURAFSKY; MARTIN, 2022). Aplicações de modelos *encoder-decoder* vão desde tradução automática à sumarização de textos. A principal ideia da arquitetura é o emprego de uma rede capaz de gerar uma representação intermediária da sequência de entrada, um contexto, que é depois utilizada para gerar uma saída apropriada para a tarefa em questão (GOODFELLOW; BENGIO; COURVILLE, 2016). Como mostrado na Figura 3.7, a sequência de entrada x_1^n é utilizada para gerar um vetor de contexto, que será utilizado para prever a saída y_1^m .

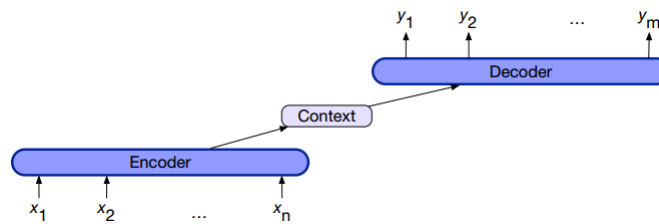


Figura 3.7: Arquitetura *Encoder-decoder*. Fonte: (JURAFSKY; MARTIN, 2022)

Essa separação entre blocos de *Encoder*, que constrói uma representação (contexto) da entrada, e *Decoder*, que utiliza o contexto criado pelo *Encoder*, simplifica a arquitetura do modelo. A Figura 3.8 exemplifica a arquitetura *Encoder-decoder* com o caso de um tradutor automático, cada palavra na língua inglesa é codificada em uma unidade do *Encoder*, cada representação é utilizada para codificar a palavra seguinte até a criação do vetor de contexto que representa toda a sentença inicial e vai ser utilizado para gerar cada palavra da sentença traduzida (JURAFSKY; MARTIN, 2022).

A arquitetura *Encoder-decoder* utilizando RNNs é retratada na Figura 3.9, em que a sentença de entrada, formada pelas palavras x_1 a x_n , é “codificada”³ no *Encoder* (utilizando uma RNN) resultando no estado final h_n^e . Esse estado final, chamado de contexto (c), será utilizado pelo *Decoder* (outra RNN) para gerar cada palavra de y_1 a y_n , junto com o estado interno anterior h_{i-1}^d e a palavra predita anteriormente y_{i-1} (CHO et al., 2014b).

Na tarefa de *Image captioning*, após a extração de uma representação da imagem, a sentença que a descreve é gerada palavra a palavra. Os autores em (VINYALS et al.,

³Durante a geração de contexto, as saídas da RNN não é utilizada.

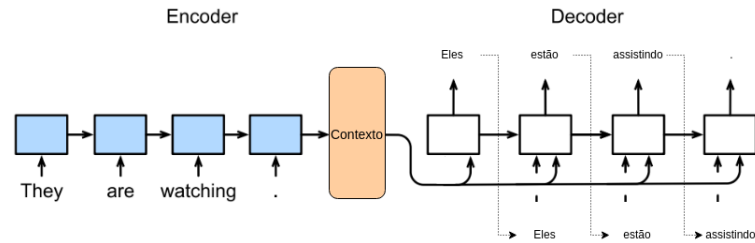


Figura 3.8: Exemplo de tradução da sentença “They are watching.” para “Eles estão assistindo.”. Fonte: adaptado de (ZHANG et al., 2021a)

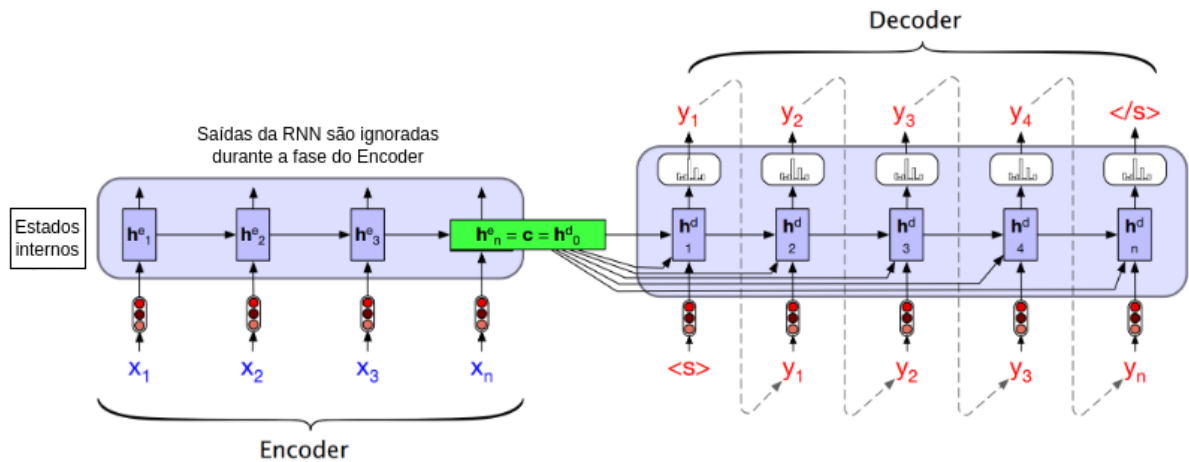


Figura 3.9: Arquitetura *Encoder-decoder*. Fonte: (JURAFSKY; MARTIN, 2022)

2015) definem a tarefa como “traduzir uma imagem”.

3.1.5 Mecanismo de Atenção

Na arquitetura *Encoder-decoder* o vetor de contexto representa uma função de todos os estados internos h_n^e do *Encoder*, conforme representado na Equação 3.1.5.1:

$$c = f(h_1^e \dots h_n^e). \quad (3.1.5.1)$$

Esse vetor de contexto armazena toda a informação sobre os n estados da entrada. Essa estratégia gera um acúmulo de informações (JURAFSKY; MARTIN, 2022). No caso de tradução automática informações do início da sentença podem não ser muito bem representadas, principalmente se for uma sentença longa (CHO et al., 2014a). Na Figura 3.9 por exemplo, as informações de h_n^e são armazenadas no vetor c , independentemente do tamanho da sequência de entrada n .

Para resolver esse problema os autores em (BAHDANAU; CHO; BENGIO, 2014) propõem que o vetor de contexto c seja adaptado a cada geração de palavra da sequência. Como visto na Figura 3.10, o vetor de contexto, agora c_1 , que será utilizado na predição

da palavra y_i pelo *Decoder* é dado por uma soma ponderada dos estados internos do *Encoder*.

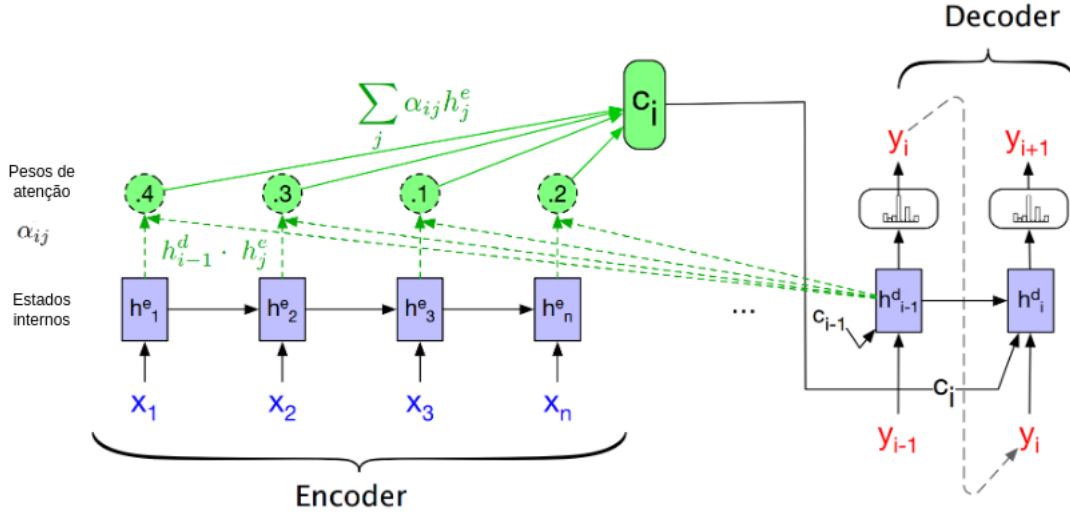


Figura 3.10: Exemplo de arquitetura *Encoder-decoder* com uso de atenção. Fonte: adaptado de (JURAFSKY; MARTIN, 2022)

Essa soma é ponderada utilizando o estado interno \mathbf{h}_{i-1}^d do *Decoder*, portanto cada vetor de contexto \mathbf{c}_i é computado levando em conta uma parte relevante da entrada que terá mais “atenção” ao prever a palavra atual. Essa relevância é mostrada na Equação 3.1.5.2.

$$\text{score}(h_{i-1}^d, h_j^e) = h_{i-1}^d \cdot h_j^e \quad (3.1.5.2)$$

A função *score* retorna o produto interno entre o estado interno da palavra anteriormente gerada pelo *Decoder* (h_{i-1}^d) e cada estado interno, gerado pelo *Encoder*, de cada palavra j da sentença de entrada (h_j^e). O uso de produto interno faz com que o mecanismo de atenção calcule uma medida de similaridade entre h_{i-1}^d e h_j^e , os autores em (BAHDANAU; CHO; BENGIO, 2014) mostram que essa similaridade concorda com a intuição humana. Para utilizar a função *score*, normaliza-se a saída com uma função *softmax*⁴ mostrada na Equação 3.1.5.3, então tem-se um vetor de pesos α_{ij} que é utilizado na soma ponderada que vai retornar o vetor de contexto c_i , como mostrado na Equação 3.1.5.4 (JURAFSKY; MARTIN, 2022).

$$\alpha_{ij} = \text{softmax}(\text{score}(h_{i-1}^d, h_j^e) \forall j \in e) \quad (3.1.5.3)$$

$$c_i = \sum \alpha_{ij} h_j^e \quad (3.1.5.4)$$

O mecanismo de atenção permite que o *Decoder* gere palavras a partir de um vetor de contexto que se adapta a cada predição anterior, focando em partes diferentes da sentença

⁴Softmax é uma função que normaliza as saídas, transformando-as em números não negativos que somados resultam em 1 (ZHANG et al., 2021a).

de entrada em cada passo de predição. O trabalho feito por (XU et al., 2016) emprega esse mesmo mecanismo para criar um vetor de contexto que se atenta a partes diferentes da imagem ao predizer cada palavra da sentença de descrição.

3.1.6 Transformers

Em (VASWANI et al., 2017), os autores retratam limitações das RNNs por conta de sua arquitetura: perda de informações devido ao uso de uma série extensa de conexões recorrentes e a dificuldade de treinamento, já que sua natureza sequencial dificulta paralelização. Após elencarem essas dificuldades, é apresentado o conceito de *Transformer*. Um *Transformer* é uma arquitetura de rede neural do tipo *Encoder-decoder* que contando inteiramente com o mecanismo de atenção para calcular as relevâncias presentes na sequência de entrada (evitando o uso de recorrência). A extração de representações de uma sentença é feita utilizando um mecanismo de atenção que relaciona entradas dessa mesma sentença em posições diferentes, por isso dá-se o nome de *Self-attention* (ou de *intra-attention*) a este mecanismo de atenção. Um *Transformer*, é capaz de gerar um vetor de contexto (também chamado de representação global), apenas com o uso de *Self-attention*, sem contar com uma RNN. Os objetivos de utilizar *Self-attention* são: reduzir a complexidade das operações feitas por RNNs, paralelizar as operações necessárias e melhorar o aprendizado de representações globais a longas distâncias. O uso de *Self-attention* permite desempenhos superiores quando comparados com RNNs que utilizam o mecanismo de atenção e com menos tempo de treinamento (VASWANI et al., 2017).

Os autores trazem uma definição diferente para o mecanismo de atenção, com o uso dos termos requisição (*query*), chave (*key*) e valor (*value*):

- requisição (*query*): é o papel dado ao foco principal em que o mecanismo está “prestando atenção” e será comparado com as entradas anteriores;
- chave (*key*): é a entrada anterior sendo comparada com a *query* atual.
- valor (*value*): é o valor dado à compatibilidade (ou relevância) de uma *query* e a *key* sendo comparada.

No contexto de *Self-attention query, key* e *value* são vetores obtidos da mesma sequência de entrada. Os autores fazem mais de uma operação de atenção com o uso de *Multi-Head Attention*, em que mecanismos de atenção diferentes são utilizados sobre projeções lineares dos vetores de *query, key* e *value* de maneira paralela (diferentes *heads*), os resultados dessas aplicações da função de atenção são concatenados e então linearmente projetados novamente. O uso de chamadas de atenção paralelas ajuda o modelo a aprender diferentes tipos de relações de relevância da entrada.

Ilustrado na Figura 3.11, um *Transformer* emprega uma arquitetura *Encoder-decoder* utilizando blocos empilhados de *self-attention* (na forma de *Multi-Head Attention*), redes neurais com conexão para frente, conexões residuais entre camadas⁵ e camadas de

⁵Conexões residuais passam informações de uma camada anterior para uma superior sem passar por uma camada intermediária (HE et al., 2015).

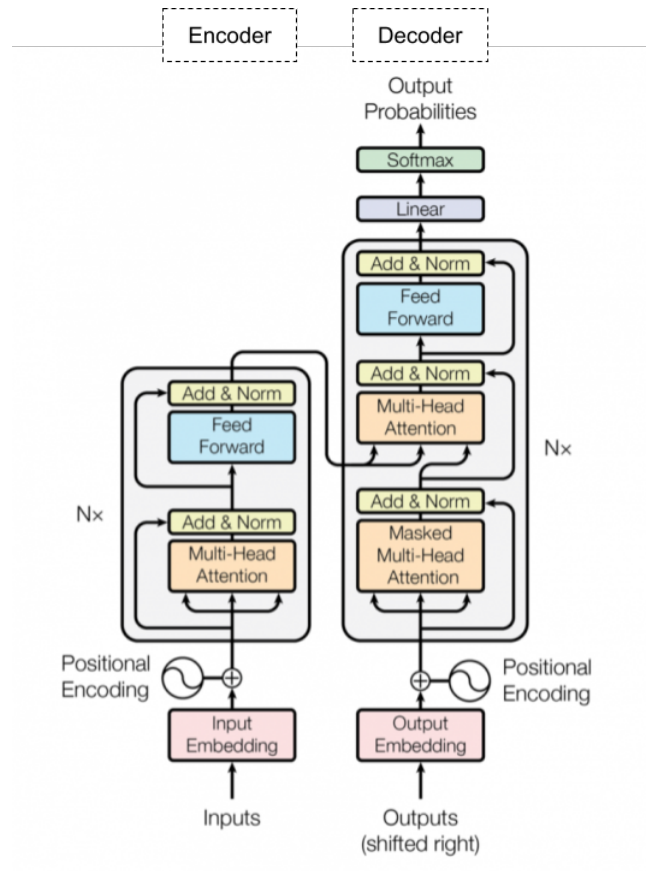


Figura 3.11: Arquitetura *Encoder-decoder* de um Transformer. Fonte: adaptado de (VASWANI et al., 2017)

normalização⁶.

A parte do *Encoder* é composta por N_x camadas empilhadas, possuindo uma camada de *Multi-Head Attention* seguida de uma camada de rede neural de conexão direta. Cada uma das camadas é seguida de uma camada de normalização e conexão residual. Essa etapa é responsável por criar dependências globais dentro da sequência de entrada inserida.

Já o *Decoder*, que também possui N_x camadas empilhadas, é composto por duas camadas de *Multi-Head Attention*. A primeira delas recebe as saídas sendo gerada pelo *Transformer*. Para evitar a criação de dependências entre a palavra atual i e as futuras entradas, duas medidas são tomadas: as sentenças de saída são deslocadas para a direita (*shifted right*), de forma que cada palavra é predita baseada nas predições anteriores apenas; além disso a primeira camada é alterada por uma máscara⁷ que impede o me-

⁶A camada de normalização é aplicada para melhorar a performance do treino, mantendo os valores dentro de um espectro não muito alto, o que facilita o aprendizado (JURAFSKY; MARTIN, 2022).

⁷Essa máscara leva os valores das dependências globais com as posições subsequentes a valores infinitamente negativos.

canismo *Self-attention* de criar dependências com posições subsequentes com as palavras que ainda serão preditas, por conta disto a camada é denominada: *Masked Multi-Head Attention*.

Após aplicar conexão residual e normalização na saída da *Masked Multi-Head Attention*, *Self-attention* é mais uma vez aplicada entre a saída do *Encoder* e as representações obtidas pela camada anterior, sendo as *queries* obtidas da camada de *Masked Multi-Head Attention*, enquanto que as *keys* e *values* vêm da saída do *Encoder*.

Em adição às camadas de *Self-attention* ambos o *Encoder* e o *Decoder* têm suas saídas dadas como entradas para duas redes neurais com conexão para frente (*Feed Forward* na imagem). No *Decoder*, à saída da rede neural com conexão para frente são aplicadas operações de normalização, transformação linear⁸ e a função *Softmax* para converter a saída para as probabilidades de cada palavra.

Como o modelo de *Transformer* não possui recorrências, é adicionada uma codificação de posição, chamada de *Positional Encoding*, nas representações de entrada e de saída. Para isso, os autores implementam funções seno e cosseno às representações das entradas das sentenças dependendo da paridade da posição de cada palavra. A hipótese é a de que os formatos das ondas ajudariam o modelo a prestar atenção em posições relativas da sentença. Com o advento dos *Transformers* o campo de *Image captioning* mudou sua atenção do uso de RNNs para o *Self-attention* (STEFANINI et al., 2021).

3.1.7 Evolução dos Métodos

Conforme o avanço das tecnologias utilizadas, as arquiteturas de *Image captioning* foram sendo adaptadas. A Figura 3.12 exemplifica os passos de extração de características de uma uma imagem por uma CNN para a geração de uma sentença de saída.

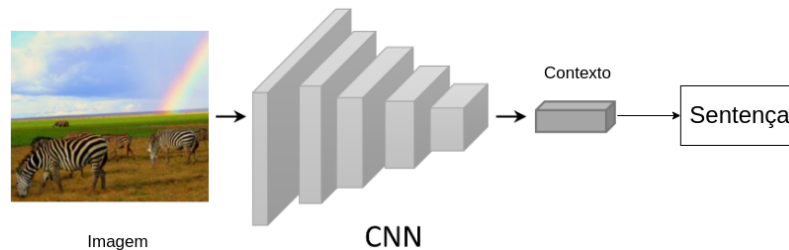


Figura 3.12: Arquitetura *Encoder-decoder* aplicada a *Image captioning* com uma CNN extraíndo representações de uma imagem e gerando o vetor de contexto. Fonte: adaptado de (STEFANINI et al., 2021)

De maneira análoga, a Figura 3.13 mostra a mesma geração, mas utilizando o mecanismo de atenção para que a geração de palavras da sentença “atente” para partes específicas da imagem.

⁸Uma transformação linear é uma correspondência que associa vetores de espaços vetoriais diferentes, preservando operações de adição vetorial e multiplicação por escalar (LIMA, 2014).

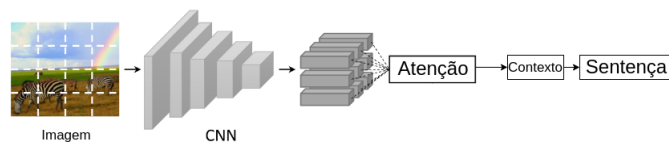


Figura 3.13: Arquitetura *Encoder-decoder* aplicada a *Image captioning* com uma CNN extraindo representações de uma imagem e utilizando atenção para gerar um vetor de contexto apropriado para cada palavra. Fonte: adaptado de (STEFANINI et al., 2021)

A Figura 3.14 ilustra o uso de *Transformers* na geração de descrições para uma imagem, utilizando *Multi-Head Attention* na saída de uma CNN, *Masked Multi-Head Attention* nas saídas geradas até a palavra atual e, mais uma vez, *Multi-Head Attention* entre as dependências globais da imagem e as palavras previstas até o passo atual.

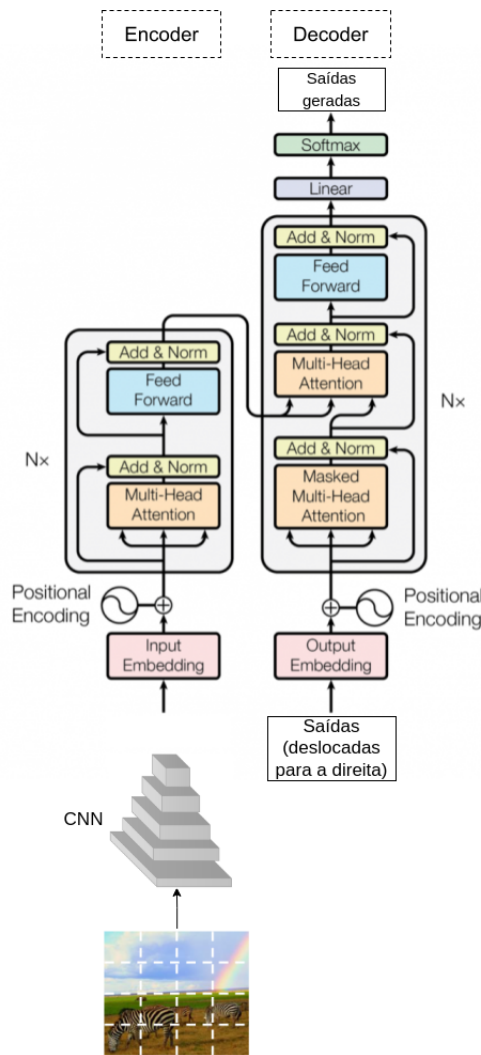


Figura 3.14: Arquitetura *Encoder-decoder* aplicada a *Image captioning* com uma CNN extraíndo representações de uma imagem e utilizando *Transformers* para criar dependências globais e gerar uma descrição para a entrada. Fonte: adaptado de (STEFANINI et al., 2021) e (VASWANI et al., 2017)

3.2 PROCESSAMENTO DE LINGUAGEM NATURAL

Processamento de Linguagem Natural (PLN) é uma área da Inteligência Artificial que estuda interações entre humanos e computadores por meio da linguagem natural (ZHANG et al., 2021a). O objetivo da pesquisa na área é compreender as diversas observações linguísticas existentes nas comunicações entre seres humanos em conversas, mensagens, escrita ou fala, caracterizando e explicando como produzimos linguagem (MANNING; SCHÜTZ, 1999).

Tarefas de PLN incluem a manipulação e entendimento automatizados de linguagem natural e podem ir desde a contagem de frequências de uma palavra no texto até a compre-

ensão da mensagem passada através da linguagem sendo analisada. PLN pode ser dividida em Geração de Linguagem Natural (do inglês *Natural Language Generation (NLG)*) e Compreensão de Linguagem Natural (do inglês *Natural Language Understanding (NLU)*) (BIRD; KLEIN; LOPER, 2009).

NLG é a área de PLN que estuda o desenvolvimento de sistemas capazes de gerar saídas de texto compreensíveis para o entendimento humano a partir de entradas que podem ser textuais ou não. Tarefas existentes em NLG são:

- Determinação de conteúdo: definir quais conteúdos incluir em um texto em construção;
- Agregação de Sentenças: definir quais informações devem aparecer em sentenças individuais;
- Realização linguística: combinar palavras para formar frases em linguagem natural e compreensível. Sendo *Image captioning* uma geração de sentença em linguagem natural a partir de uma entrada de imagem, a tarefa é um paradigma que se encaixa dentro de NLG (GATT; KRAHMER, 2018).

NLU é a área de PLN que pesquisa o entendimento do texto dado como entrada, representando a entrada em informações formais e estruturadas (TSENG et al., 2020). NLU busca compreender a estrutura do texto, respondendo questões como “o que aconteceu?” ou “quem fez o que para quem?” (BIRD; KLEIN; LOPER, 2009). Exemplos de tarefas de NLU de acordo com autores em (WANG et al., 2018) são:

- Análise de Sentimento: dada uma sentença, classificar o que está escrito em positivo ou negativo;
- Similaridade semântica de sentenças: verificar se dois pares de sentenças são equivalente semanticamente;
- Tarefas de inferência e extração de informação: dados dois fragmentos de texto, inferir quais as relações entre eles;
- Resposta de perguntas: dada uma sentença, responder um questionamento feito sobre ela.

A Figura 3.15 ilustra as possíveis análises linguísticas que se pode fazer de uma sentença. Neste trabalho os níveis morfológico e sintático serão analisados.

3.2.1 Análise Morfológica

Morfologia pode ser compreendida como o estudo da forma das palavras bem como suas estruturas internas (BENDER, 2013). O estudo da morfologia dentro do campo de PLN é importante pois, em diferentes conjuntos de dados, é comum deparar-se com exemplos de palavras que ainda não foram vistos anteriormente, mas que podem estar morfológicamente associadas a outras palavras conhecidas (vistas previamente no processo de

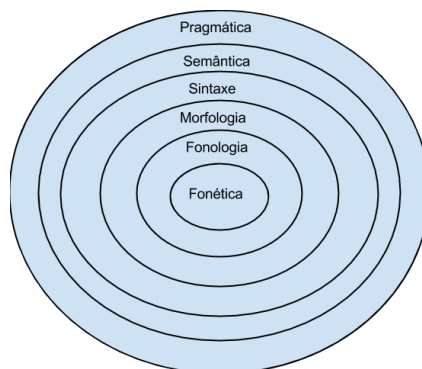


Figura 3.15: Diagrama com as análises que podem ser feitas de uma sentença. Fonte: (SILVA, 2010)

treinamento de um modelo de geração de texto por exemplo) (MANNING; SCHÜTZE, 1999). Portanto, ao entender processos morfológicos, um modelo de PLN pode inferir conhecimento sobre a sintaxe ou o sentido de uma palavra (MANNING; SCHÜTZE, 1999).

Dentro da morfologia, um conceito importante para análises é o do morfema: uma unidade mínima que possui significado, não podendo ser dividida sem que seu significado seja alterado (SILVA, 2010). Palavras monomorfêmicas são constituídas por um único morfema (por exemplo “mãe” ou “feliz”), enquanto que palavras constituídas por mais de um morfema (como “mães” ou “infelizmente”) são chamadas polimorfêmicas (SILVA, 2010).

Alguns atributos específicos são usados para classificar morfemas em diferentes grupos. Ao analisar a palavra “ruas”, formada pelos morfemas “rua” e “-s”, percebe-se que o morfema “rua”, por si só, também pode ser utilizado como uma palavra, enquanto que o morfema “-s”, que representa a forma plural em “ruas”, não pode figurar como palavra isoladamente (CUNHA; CINTRA, 2016). Morfemas que podem aparecer sozinhos como palavras são chamados morfemas livres, enquanto que morfemas que nunca são utilizados de forma isolada são chamados morfemas presos (CUNHA; CINTRA, 2016).

Há também uma designação de morfemas baseada em seus significados. Os morfemas lexicais possuem significado externo ao mundo linguístico: pessoas, instrumentos, seres, ideias (SILVA, 2010). Esse tipo de morfema forma uma classe aberta, ou seja, passível de ser aumentada com a adição de novos componentes (CUNHA; CINTRA, 2016). Já os morfemas gramaticais, que derivam das relações linguísticas existentes em um idioma, têm significação interna e pertencem a um paradigma fechado, onde o número de elementos do grupo é restrito no idioma (SILVA, 2010). Uma consequência dos grupos aberto e fechado dos morfemas é a maior frequência de morfemas gramaticais ao se examinar um texto, diferente do que é observado com morfemas lexicais que têm frequência menor (CUNHA; CINTRA, 2016).

Com relação a construção de palavras, aos morfemas lexicais geralmente são dados o nome de “radical”. A partir de um radical, palavras de uma mesma família são unidas

com um significado base compartilhado (CUNHA; CINTRA, 2016). Ao radical pode-se adicionar os morfemas gramaticais, que podem aparecer de diferentes formas: desinência (também chamada de morfema flexional), afixos (ou morfemas derivacionais) ou como vogal temática (CUNHA; CINTRA, 2016).

Os morfemas flexionais (desinências) indicam: gênero e número de substantivos, adjetivos e alguns pronomes, o número e a pessoa de um verbo. Existe, portanto, desinências nominais e desinências verbais na língua portuguesa, sendo as nominais de gênero masculino (“-o”) e feminino (“-a”) ou de número: plural (“-s”) e singular que é expressa pela “desinência-zero” ou a falta de desinência (CUNHA; CINTRA, 2016). Por exemplo: nas palavras “frias” e “fazemos”, identifica-se a presença das desinências “-a” (caracterizando a forma feminina em “frias”), “-s” (representado a forma plural em “frias”) e “-mos” (indicando a 1ª pessoa do plural em “fazemos”) (CUNHA; CINTRA, 2016).

As desinências exemplificadas acima não alteram a classe gramatical das palavras, o que é uma característica de morfemas flexionais em contraste com os morfemas derivacionais (SILVA, 2010). Os morfemas derivacionais, também chamados de “afixos”, modificam o sentido do radical ao qual são adicionados e são chamados de “prefixos” ou “sufixos” a depender de sua posição com relação ao radical, respectivamente: antes ou depois do mesmo (CUNHA; CINTRA, 2016).

Por fim, a vogal temática indica a qual das três conjugações existentes no português um verbo pertence. No exemplo de “fazemos”, identifica-se os morfemas “faz-” e “mos”, havendo ainda a presença de um outro morfema⁹: a vogal temática “-e-”, que caracteriza o verbo como pertencente à segunda conjugação. Tem-se, portanto, as seguintes vogais temáticas e suas respectivas conjugações: “-a-” (1ª conjugação), “-e-” (2ª conjugação) e “-i-” (3ª conjugação) (CUNHA; CINTRA, 2016).

3.2.2 Análise Sintática

A análise sintática consiste em atribuir uma estrutura sintática a uma sentença, sendo essa estrutura definida pelas classes gramaticais dos termos desta sentença. Analisar sintaticamente uma sentença é um passo intermediário para tarefas de PLN como: verificação gramatical, tradução automática, análise de dependência ou responder perguntas sobre uma sentença (JURAFSKY; MARTIN, 2022). A análise estrutural dos termos sintáticos de uma sentença pode ser feita derivando a sequência de palavras que a constitui. Assim, pode-se representar o resultado de uma análise de classes gramaticais em formato de árvore enraizada¹⁰, com as classes gramaticais definidas pela análise (MANNING; SCHÜTZE, 1999), um exemplo de árvore sintática pode ser visto na Figura 3.16, onde à frase “Eu prefiro comer feijão com arroz” são atribuídas as classes gramaticais de cada palavra.

⁹Os autores em (CUNHA; CINTRA, 2016) observam que não há um consenso quanto a inclusão de vogais temáticas entre morfemas. Como, neste trabalho, as análises foram feitas observando-se morfemas flexionais, a inclusão ou não de vogais temáticas como morfemas não será aprofundada.

¹⁰Uma árvore enraizada é uma representação de um gráfico acíclico em que um dos vértices é destacado (denominado de raiz) (CORMEN et al., 2009)

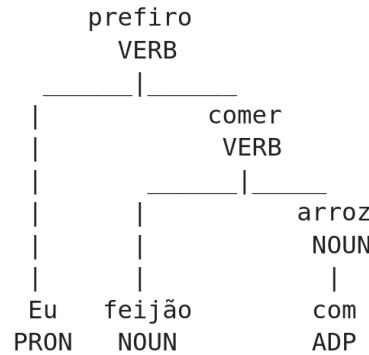


Figura 3.16: Árvore sintática para a frase “Eu prefiro comer feijão com arroz”. Fonte: preparado pelos autores.

3.2.2.1 Classes Gramaticais

Classes gramaticais (*Parts of Speech* - Part-of-speech (POS)) são agrupamentos de palavras que possuem similaridade de função sintática¹¹ entre elas (MANNING; SCHÜTZ, 1999). Alguns exemplos de POS são: nome, advérbio, verbo, adjetivo, determinante. Algumas palavras podem possuir mais de uma classe gramatical a depender do seu uso na sentença, por exemplo nas sentenças:

- “fez o **bem** sem olhar a quem”
- “os trabalhos ficaram **bem** feitos”

A palavra “bem” pertence à classe gramatical **substantivo** na primeira frase e **advérbio** na segunda.

Classes gramaticais podem ser divididas em dois grupos: variáveis e invariáveis (CUNHA; CINTRA, 2016). Essa divisão ocorre se uma classe gramatical permite ou não a combinação com morfemas.

As classes gramaticais (POS) utilizadas neste exemplo, e na realização do trabalho, são as definidas em (RADEMAKER et al., 2017). A escolha foi feita para facilitar a reprodutibilidade, uma vez que a extração de classes gramaticais é feita com o *framework spaCy* (MONTANI et al., 2022) que utiliza a mesma fonte. *spaCy*¹² (MONTANI et al., 2022) é uma biblioteca de ferramentas desenvolvida para PLN. As ferramentas do *spaCy* permitem desenvolvimento e compreensão de dados textuais, podendo ser utilizadas para diferentes tarefas da área como: extração de informação, NLU ou pré-processamento de texto com aprendizado profundo. Neste trabalho, *spaCy* foi utilizado para extração de informações de classes gramaticais que serão utilizadas para treinar o modelo de *Image captioning*.

A Tabela 3.2 mostra os rótulos das classes gramaticais utilizadas e seus respectivos significados.

¹¹A sintaxe estuda as combinações e relações entre termos de uma oração (BECHARA, 2009).

¹²<https://spacy.io/>

Rótulo	Significado
ADJ	Adjetivo: palavras que modificam nomes e especificam suas propriedades ou atributos. Exemplo: verde, grande, bonito.
ADP	Adposição: abrange a classe das preposições, posposições e circumposições. Exemplo: com, de, em.
ADV	Advérbio: palavras que modificam verbos. Exemplo: muito, bem, exatamente.
AUX	Verbo auxiliar: verbo que acompanha o verbo principal de uma frase. Exemplo: estou lendo, foi feito
CONJ	Conjunção coordenativa: palavra que liga palavras expressando uma relação semântica. Exemplo: e, ou, mas
DET	Determinante: palavras que modificam substantivos (artigos, determinantes possessivos, determinantes demonstrativos, determinantes quantificadores). Exemplo: meu, um, aquele.
INTJ	Interjeição: palavra normalmente usada como exclamação ou parte de uma exclamação. Exemplo: ei!, opa!
NOUN	Substantivo: classes gramaticais que denotam uma pessoa, lugar, coisa, animal ou ideia. Exemplo: garota, cachorro, feijão.
NUM	Numeral: palavra funcionando como determinante, adjetivo ou pronome, que expressa um número e a relação com o número, como quantidade, sequência, frequência ou fração. Exemplo: um, 2, III.
PRON	Pronomes: palavras que substituem o uso de substantivos, cujo significado pode ser inferido pelo contexto linguístico da frase. Exemplo: esta, mim, isto.
PROPN	Nome próprio: substantivo que é o nome de um indivíduo específico, lugar ou objeto. Exemplo: Maria, Salvador, Pituba.
PUNCT	Pontuação: marca caracteres que delimitam o texto. Exemplo: , . ;
SCONJ	Conjunção subordinativa: palavras cujo função é unir orações. Exemplo: que, embora.
VERB	Verbo: classe de palavras que sinalizam eventos e ações. Exemplo: correr, comer, andar.
X	Outros.

Tabela 3.2: Tabela com os rótulos e significados de cada POS de acordo com (RADEMAKER et al., 2017).

3.3 MÉTRICAS DE AVALIAÇÃO

As métricas de avaliação específicas para a tarefa de *Image Captioning* surgem apenas a partir de 2014 com a *Consensus-based Image Description Evaluation (CIDEr)* (VEDANTAM; ZITNICK; PARIKH, 2014), seguida da *Semantic Propositional Image Caption Evaluation (SPICE)* (ANDERSON et al., 2016). Antes disso, a avaliação de modelos de descrição de imagem era feita utilizando ferramentas idealizadas para a tarefa de tradução automática: Bilingual Evaluation Understudy (BLEU) e Metric for Evaluation of Translation with Explicit Ordering (METEOR) (BERNARDI et al., 2017). Essas duas métricas ainda são utilizadas para avaliação comparativa de modelos e sendo assim utilizadas neste

trabalho.

Métricas de avaliação de NLG por vezes utilizam o *n-gram*. Um *n-gram* é uma sequência de n palavras, ou seja, um *1-gram* (também chamado de *unigram*) é uma sequência de uma palavra (“eu”), um *2-gram* (*bigram*) é uma sequência de duas palavras (“eu prefiro”), um *3-gram* (*trigram*) corresponde a uma sequência de três palavras (“eu prefiro comer”) e assim por diante. *N-grams* podem ser utilizados como modelos de linguagem que calculam a probabilidade da última palavra de uma sequência dadas as palavras anteriores (JURAFSKY; MARTIN, 2022). As métricas de avaliação que aparecem nesta seção utilizam até *4-grams* para medir o desempenho de um modelo fazendo comparações entre sentenças geradas e sentenças de referência.

3.3.1 BiLingual Evaluation Understudy - BLEU

Bilingual Evaluation Understudy (BLEU) (PAPINENI et al., 2002) é uma métrica de qualidade primeiramente introduzida para a tarefa de tradução de máquina (*machine translation*) e é baseada na contagem de precisão (*precision*) de *n-grams* (a contagem pode ser de 1, 2, 3 ou 4 *n-grams*) entre a sentença gerada (sentença candidata) e as sentenças de referência (como normalmente há mais de uma tradução ideal para uma sentença, a depender da escolha de palavras ou ordem utilizada, usa-se mais de uma sentença de referência para avaliar uma tradução (PAPINENI et al., 2002)).

A métrica BLEU pressupõe que se uma sentença candidata compartilha muitas palavras com uma sentença referência, ela é uma candidata de qualidade, portanto a métrica tenta parear *n-grams* entre duas sentenças e contar a quantidade de combinações encontradas. Quanto mais combinações, melhor é a qualidade da tradução (PAPINENI et al., 2002).

Porém, sistemas de tradução automática podem gerar sentenças candidatas que possuem palavras repetidas em comum com uma sentença referência, a contagem de *n-grams* feita pela métrica BLEU é modificada para “exaurir” um *n-gram* já contado. A Equação 3.3.1.1 mostra como isso é feito: primeiro calcula-se a quantidade combinações de *n-grams* presentes na sentença candidata e nas sentenças referências, a contagem de candidatos é cortada (*clipped*) pela quantidade máxima de *n-grams* presentes nas sentenças referências¹³, então divide-se pelo total de *n-grams* da sentença candidata, sendo este passo chamado de “precisão modificada”(PAPINENI et al., 2002).

$$p_n = \frac{\sum_{C \in \{candidatas\}} \sum_{n\text{-grams} \in C} Count_{clip}(n\text{-grams})}{\sum_{C' \in \{candidatas\}} \sum_{n\text{-grams} \in C'} Count(n\text{-grams})} \quad (3.3.1.1)$$

Uma sentença traduzida não deve ser nem muito maior e nem muito menor que as sentenças referência. A contagem de *n-grams* penaliza sentenças muito longas que não possuem palavras em comum com a referência e a precisão modificada também é penalizada se uma sentença candidata repete muitas palavras compartilhadas com a referência. No entanto, caso uma sentença candidata seja muito curta, a precisão modificada falha em avaliar corretamente o comprimento da sentença e a cobertura (*recall*) das palavras,

¹³ $Count_{clip} = \min(Count, Max_Ref_Count)$, neste passo a contagem de palavras é limitada para não ultrapassar a contagem máxima de palavras em alguma das sentenças de referência.

por conta disto é adicionada uma “penalidade por brevidade” (do inglês *Brevity Penalty* - BP). Por conta desta penalidade por brevidade, traduções candidatas devem ter tamanho, escolha e ordenação de palavras que combinem com a sentença de referência (PAPINENI et al., 2002). A Equação 3.3.1.2 mostra o cálculo de BP, onde r e c são a quantidade de palavras no texto de referência e a quantidade de palavras no texto candidato respectivamente:

$$BP = \begin{cases} 1 & \text{se } c > r \\ e^{\frac{1-r}{c}} & \text{se } c \leq r \end{cases} \quad (3.3.1.2)$$

A métrica BLEU é então calculada pelo produto da média geométrica ($w_n = 1/N$) da precisão modificada com o fator de penalidade por brevidade, como mostrado na Equação 3.3.1.3.

$$BLEU = BP * exp\left(\sum_{n=1}^N w_n \log p_n\right) \quad (3.3.1.3)$$

O uso de diferentes tamanhos de n -grams na precisão modificada se deve a tentativa de capturar aspectos diferentes da tradução: adequação e fluência. Uma tradução que usa as mesmas palavras (1 -gram) que a sentença de referência é adequada, já uma sentença que combine n -grams mais longos leva em consideração a fluência.

Para exemplificar a métrica BLEU, usando as seguintes sentenças:

- sentença a ser traduzida: “*she read the book because she was interested in world history*”
- sentença referência: “ela leu o livro porque ela estava interessada na história mundial”
- sentença candidata 1: “ela leu o livro porque ela estava interessada na história mundial”
- sentença candidata 2: “ela estava interessada na história mundial porque ela leu o livro”

Considerando que a sentença candidata é uma tradução perfeita, já que é igual a sentença referência 1, a métrica BLEU para 1 -grams, 2 -grams, 3 -grams e 4 -grams é a nota máxima 1. Já para a sentença candidata 2, que utiliza as mesmas palavras da referência porém em uma ordem diferente, as pontuações são: 1 para 1 -gram, 0.9 para 2 -gram, 0.66 para 3 -gram e -0.5 para 4 -gram, nota-se que, quanto maior a sequência de palavras que se espera que combinem entre candidata e referência, pior a nota dada pela métrica BLEU.

3.3.2 Metric for Evaluation of Translation with Explicit Odering - METEOR

Metric for Evaluation of Translation with Explicit Odering (METEOR) (BANERJEE; LAVIE, 2005) é uma métrica de tradução automática feita para superar problemas detectados na métrica BLEU. Ela supera a métrica BLEU em julgamento humano pois, de acordo com os autores, favorece o *recall* na construção da métrica. *Recall* (a quantidade de combinações de *n-grams* do total de *n-grams* presentes na sentença de referência) é importante para avaliar a saída de uma tradução pois reflete a cobertura da tradução feita com relação a referência, na métrica BLEU essa cobertura é compensada com a penalidade por brevidade, mas ela não é suficiente para compensar a falta de cálculo do *recall*. METEOR é desenhada para resolver essa fraqueza, utilizando combinações palavra a palavra (*1-gram*) entre tradução e sentença referência.

Dada uma sentença candidata e uma sentença de referência, na primeira fase, as palavras de cada sentença são alinhadas (um alinhamento é definido como um mapeamento de *unigram*, em que cada *unigram* de uma sentença é mapeado para zero ou um *unigram* de outra). Esse mapeamento é repetido em estágios com módulos diferentes mapeando *unigrams* de acordo com um critério específico em cada uma dessas fases: o primeiro estágio mapeia o exato *unigrams*, o segundo mapeia *unigrams* sinônimos entre si, outro estágio mapeia *unigrams* se eles tiverem a mesma derivação¹⁴ (*stemming*). Como um mapeamento de *unigrams* só acontece se um *unigrams* não tiver sido mapeado, a ordem em que os módulos são utilizados nos estágios impõe uma prioridade diferente para cada módulo, ou seja, se o primeiro módulo mapear *unigrams* “exatos”, a métrica irá dar uma nota maior para traduções que usem as mesmas palavras da sentença de referência.

A segunda fase recolhe o maior subconjunto de mapeamento de *unigrams* resultante da fase anterior e calcula uma média harmônica de acordo com a Equação 3.3.2.1, em que P é a fração de *unigrams* da sentença candidata que aparecem na sentença referência, R é a fração de *unigrams* da sentença referência que aparecem na sentença candidata e α é um parâmetro que define a qual fração será dada mais importância.

$$Fmean = \frac{P * R}{\alpha * P + (1 - \alpha) * R} \quad (3.3.2.1)$$

Como a Equação acima leva em conta apenas valores calculados com *unigrams*, METEOR também calcula uma penalidade para considerar combinações mais extensas. Essa penalidade é mostrada na Equação 3.3.2.2, em que a quantidade de fragmentos (grupos de *unigrams* adjacentes da sentença candidata que foram mapeados para *unigrams* também adjacentes na sentença referência) é dividida pelo total de *unigrams* mapeados.

$$Penalidade = 0.5 * \left(\frac{\text{quantidade de fragmentos}}{\text{quantidade de unigrams mapeados}} \right)^3 \quad (3.3.2.2)$$

Por fim a métrica é calculada como mostrado na Equação 3.3.2.3

$$METEOR = Fmean * (1 - Penalty) \quad (3.3.2.3)$$

¹⁴Palavra após remoção do sufixo. Exemplo: felizmente → felizmente (JURAFSKY; MARTIN, 2022).

Comparando a métrica BLEU com a METEOR e com uma pequena alteração nas sentenças percebe-se que a METEOR pode corresponder melhor com o julgamento humano. Dadas as sentenças:

- sentença a ser traduzida: “*she read the book because she was interested in world history*”
- sentença referência: “ela leu o livro porque ela estava interessada na história mundial”
- sentença candidata 1: “ela leu o livro porque ela estava interessada na história do mundo”

Enquanto o BLEU para 4 -grams é de 0.79, a métrica METEOR é de 0.92, indicando que a sentença candidata 1 é uma boa tradução.

3.3.3 Consensus-based Image Description Evaluation - CIDEr

Com o avanço nas pesquisas na área de *Image captioning*, métricas especificamente criadas para avaliar uma descrição de imagem gerada automaticamente foram criadas. Uma delas é a CIDEr (*Consensus-based Image Description Evaluation*) (VEDANTAM; ZITNICK; PARIKH, 2014), uma métrica que calcula a similaridade por cosseno entre n -grams ponderados.

O objetivo da métrica CIDEr é calcular, para uma imagem I_i , qual o consenso entre a descrição candidata c_i e um conjunto de descrições $S_i = \{S_{i1}, \dots, S_{im}\}$. Quanto maior esse consenso, melhor a descrição. Para mensurar o consenso entre sentenças, todas as palavras são mapeadas para as suas formas derivadas (*stemming*). Segundo os autores, um “consenso” pode ser medido de acordo com quantidade de n -grams em comum entre c_i (sentença candidata) e S_i (sentenças referência). Além disso, n -grams não presentes em S_i também não deveriam estar em c_i , porém n -grams que aparecem constantemente em descrições de imagens (que sejam diferentes no conjunto de imagens), devem ter um peso menor, já que esses n -grams podem ser pouco informativos sobre a imagem (n -grams que são visualmente irrelevantes). Para diminuir o peso dado a n -grams pouco informativos e que aparecem muito no conjunto de sentenças de imagens diferentes, é utilizado *Term Frequency Inverse Document Frequency* (TF-IDF) (ROBERTSON, 2004) para ponderar cada n -grams. TF-IDF é calculado, para cada n -grams w_k que ocorre em uma sentença referência S_{ij} ou candidata c_i e é representado respectivamente por $h_k(S_{ij})$ e $h_k(c_i)$, sendo esse cálculo feito pela Equação 3.3.3.1, onde ω é o vocabulário de todos os n -grams e I o conjunto de todas as imagens do conjunto de dados. Nesta equação, o primeiro termo (TF) aplica um peso maior a n -grams que frequentemente aparecem na sentença dada como referência da descrição da imagem, enquanto que o segundo (IDF) aplica pesos menores a n -grams que aparecem frequentemente em todas as descrições de todas as imagens do conjunto de dados.

$$g_k(S_{ij}) = \frac{h_k(S_{ij})}{\sum_{w_l \in \Omega} h_l(S_{ij})} \log \left(\frac{|I|}{\sum_{I_p \in I} \min(1, \sum_q h_k(S_{pq}))} \right) \quad (3.3.3.1)$$

A métrica CIDEr para n -grams de tamanho n é então calculada pela Equação 3.3.3.2 que representa a média da similaridade cosseno (produto interno) entre a sentença candidata e as sentenças referências, o que leva em conta tanto precisão quanto cobertura (*recall*), sendo $g^n(c_i)$ o vetor formado por $g_k(S_{ij})$ para todos os n -grams de tamanho n e $\|g^n(c_i)\|$ com a magnitude do vetor $g^n(c_i)$, o mesmo aplicado a $g^n(S_{ij})$.

$$CIDEr_n(c_i, S_i) = \frac{1}{m} \sum_j \frac{g^n(c_i) \cdot g^n(S_{ij})}{\|g^n(c_i)\| \|g^n(S_{ij})\|} \quad (3.3.3.2)$$

Usa-se n -grams de ordem maior (maior n) para capturar melhores semânticas e propriedades gramaticais, combinando as métricas para n -grams variados conforme:

$$CIDEr_n(c_i, S_i) = \sum_{n=1}^N w_n CIDEr_n(c_i, S_i) \quad (3.3.3.3)$$

Com pesos uniformes $w_n = 1/N$, a métrica CIDEr, ao utilizar consenso entre anotadores, captura melhor a intuição humana ao descrever uma imagem, além de se concentrar em palavras que são mais informativas para a imagem cuja sentença está sendo avaliada, dando menos importância se palavras muito comuns em todo o conjunto de dados aparecem ou não.

3.3.4 Semantic Propositional Image Caption Evaluation - SPICE

Os autores em (ANDERSON et al., 2016) apresentam a hipótese de que o componente “conteúdo semântico” é importante para a avaliação de *Image captioning* e propõem o SPICE (Semantic Propositional Image Caption Evaluation), que utiliza grafos de cena como métrica de avaliação. Com o uso de grafos de cena, a métrica SPICE é capaz de responder conteúdos específicos de modelos sendo avaliados como: “um modelo é capaz de contar?” ou “qual modelo enxerga melhor as cores da imagem?”.



Figura 3.17: Imagem de exemplo para geração de grafo de cena. Fonte: (ANDERSON et al., 2016).

Dada uma sentença candidata c e um conjunto de sentenças de referência $S = S_1, \dots, S_n$ para a Figura 3.17, primeiro uma representação intermediária das sentenças é gerada, capaz de codificar conteúdo semântico de cada uma delas. A representação escolhida pelos autores são os grafos de cena (*scene graphs*), uma representação estruturada

que é capaz de expressar objetos, atributos e relacionamentos entre os objetos e a cena (CHANG et al., 2021).

Os grafos de cena de uma sentença candidata e um conjunto de sentenças referência são, respectivamente, denotados por $G(c)$ e $G(S_i)$ para cada $S_i \in S$, os grafos de cena $G(S_i)$ são combinados em nós sinônimos. Uma sentença candidata é representada como grafo de cena pela Equação 3.3.4, onde $O(c)$, $E(c)$ e $K(c)$ são respectivamente, o conjunto de objetos mencionados em c , o conjunto de arestas que representam as relações entre objetos e o conjunto de atributos associados a esses objetos.

$$G(c) = \langle O(c), E(c), K(c) \rangle \quad (3.3.4.1)$$

Após transformar a sentença c em um grafo de cena, uma função T é definida para retornar tuplas lógicas dos grafos gerados:

$$T(G(c)) \triangleq O(c) \cup E(c) \cup K(c) \quad (3.3.4.2)$$

Cada tupla contém os elementos que representam os objetos que compõem a cena, seus atributos e relações com outros objetos da cena, por exemplo, a Figura 3.17, é representada pelas seguintes tuplas: (garota), (quadra), (garota, jovem), (garota, em pé), (quadra, tênis), (garota, em cima de, quadra).

A operação \otimes , definida como uma função que retorna tuplas coincidentes, é usada para calcular precisão e *recall* em:

$$P(c, S) = \frac{|T(G(c)) \otimes T(G(S))|}{|T(G(c))|} \quad (3.3.4.3) \quad R(c, S) = \frac{|T(G(c)) \otimes T(G(S))|}{|T(G(S))|} \quad (3.3.4.4)$$

Em seguida a métrica SPICE é calculada como a média harmônica em:

$$SPICE(c, S) = \frac{2 * P(c, S) * R(c, S)}{P(c, S) + R(c, S)} \quad (3.3.4.5)$$

Para tuplas coincidentes são utilizados sinônimos e *stemmings* em uma maneira semelhante à feita na métrica METEOR. Logo, tuplas são consideradas semelhantes se suas palavras derivadas (*stemmings*) forem iguais ou se essas palavras forem sinônimos.

3.4 RESUMO DO CAPÍTULO

Este capítulo apresenta os conceitos básicos para entendimento de um sistema de *Image captioning*. Como este trabalho tem como objetivo construir uma arquitetura para a tarefa de descrever uma imagem, as Redes Neurais (3.1) foram introduzidas e métricas de avaliação de *Image captioning* foram apresentadas. Além disso, algumas informações linguísticas são apresentadas que serão utilizadas para guiar o treinamento da arquitetura. No próximo capítulo, o estado da arte de *Image captioning* é introduzido, começando com os primeiros sistemas de regras e preenchimento de *templates*, em seguida aprendizado profundo e *Transformers* são apresentados.

TRABALHOS RELACIONADOS

Este capítulo apresenta a evolução dos sistemas de *Image captioning* até o atual estado da arte. No início, são exploradas as primeiras estratégias utilizadas para descrever uma imagem (com preenchimentos de *templates* e recuperação de sentenças de imagens similares) e seguimos para o uso mais atual de aprendizado profundo, separando os trabalhos entre arquiteturas iniciais e mais recentes com uso de *Transformers* para geração de sentenças. Em seguida, os conjuntos de dados existentes para a tarefa de descrição de imagens são apresentados.

4.1 IMAGE CAPTIONING

Esta seção apresenta a evolução dos sistemas de *Image captioning* até o atual estado da arte. Inicialmente exemplificando as primeiras estratégias utilizadas para descrever uma imagem (com preenchimentos de *templates* e recuperação de sentenças de imagens similares), em seguida para o uso mais atual de aprendizado profundo (do inglês *deep learning*), separando os trabalhos entre as arquiteturas iniciais, que aplicam o modelo *Encoder-decoder* com uso de CNNs (para extrair características da imagem de entrada) e RNNs (como modelos de linguagem), as arquiteturas que usam atenção entre CNNs e RNNs dando ao sistema a capacidade de “focar” em uma parte da imagem ao inferir a próxima palavra e os sistemas mais recentes com uso de *Transformers* para criação de dependências globais e geração de sentenças a partir de características extraídas de uma CNN.

4.1.1 Abordagens Iniciais

Os primeiros trabalhos da área se baseiam em regras e preenchimento de *templates* (YAO et al., 2010; AKER; GAIZAUSKAS, 2010; YANG et al., 2011). Esses sistemas, buscam substituir palavras em frases, buscando substantivos e suas dependências de acordo com a probabilidade máxima analisada dentre as sentenças do conjunto de dados usado para treino.

Outra abordagem é a de recuperação de descrições baseada em espaços multimodais treinados com conjuntos de dados de pares imagens-texto (KARPATHY; JOULIN; FEI-FEI, 2014; ORDONEZ; KULKARNI; BERG, 2011; FARHADI et al., 2010; PAN et al., 2004). Nessa abordagem, a similaridade de diferentes características da imagem é utilizada para buscar sentenças de imagens semelhantes e descrever a imagem dada como entrada.

Atualmente, os modelos de descrição de imagem são baseados em redes neurais profundas (STEFANINI et al., 2021).

4.1.2 Image captioning com Redes Neurais

O uso de redes neurais profundas diferencia as formas de empregar as representações extraídas de imagens na geração de sentenças. Os modelos utilizam a arquitetura *Encoder-decoder*, na qual as entradas são imagens representadas por vetores de *pixels*, que passam pelo *Encoder* e geram uma representação intermediária (*Embedding*) que é utilizada pelo modelo de linguagem do *Decoder* para gerar a descrição (BERNARDI et al., 2017). Os primeiros sistemas utilizavam uma CNN como *Encoder* e uma RNN como *Decoder*, passando o vetor de representação diretamente como entrada para o modelo de linguagem.

4.1.2.1 CNN → RNN

Em (VINYALS et al., 2015) os autores treinam um modelo para maximizar a probabilidade de uma sentença alvo dada uma imagem de entrada. No artigo, são empregadas duas redes neurais profundas: uma Rede Neural Convolutiva (*Convolutional Neural Network* - CNN) para extrair características (*features*) de imagens. Essas características são utilizadas para produzir a sentença final, palavra a palavra, com uma Rede Neural Recorrente (*Recurrent Neural Network* - RNN).

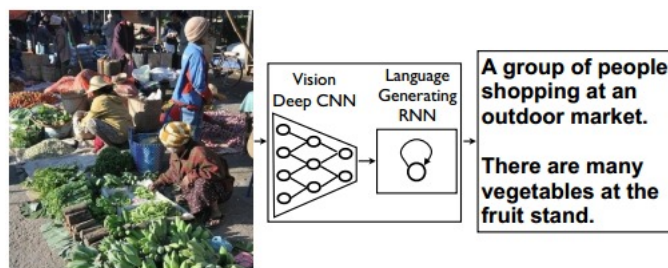


Figura 4.1: Exemplo de descrição de imagem com CNN e RNN em (VINYALS et al., 2015)

A Figura 4.1 demonstra um exemplo de geração de descrição para uma imagem, em que as características extraídas no passo do *Encoder* são utilizadas em cada passo de predição de palavra feita pela RNN. Ainda que utilize redes neurais profundas para codificar a imagem e gerar uma sentença, o trabalho utiliza a mesma representação gerada

da imagem em cada predição de palavra, o que é prejudicial por não adaptar a informação de acordo com a relevância de cada parte da imagem para a predição atual.

Autores em (HE et al., 2017), utilizam a mesma arquitetura adicionando informações de classes gramaticais das palavras, tratando as classes gramaticais das palavras como uma “pista” de informação para quando o *Decoder*, por meio da RNN, deve utilizar informações visuais para prever a próxima palavra. Como mostrado na Figura 4.2, a informação da imagem extraída pela CNN é inserida na RNN apenas em alguns passos de predição. A decisão de inserir informação visual é tomada de acordo com o tipo de classe gramatical da palavra predita anteriormente.

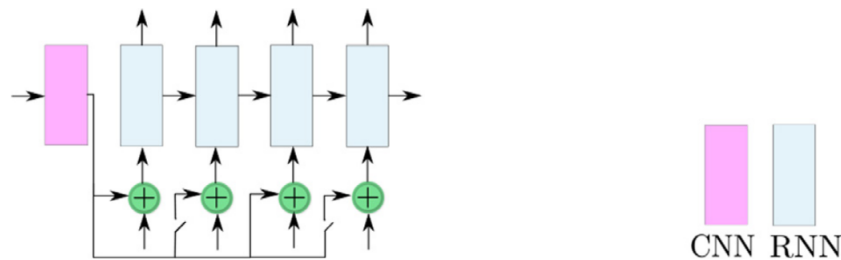


Figura 4.2: Arquitetura do modelo em (HE et al., 2017). Fonte: (HE et al., 2017)

Diferente deste trabalho, a arquitetura proposta não utiliza o mecanismo de atenção em conjunto com as características da imagem. O uso de atenção, em conjunto com a informação gramatical da palavra anterior, permite que o modelo de linguagem utilize informações visuais focadas em uma parte da imagem relevante para a palavra atual, de acordo com a POS da palavra anterior.

Em (ZHANG et al., 2021b), os autores utilizam classes gramaticais de palavras para auxiliar que o modelo treinado aprenda padrões das sequências de informações gramaticais inseridas. Uma RNN é encarregada de fazer a predição sintática da próxima palavra que será predita por outra rede RNN, à qual são dados como entrada as informações visuais, extraídas com uso de uma CNN, a informação do modelo de linguagem com as palavras preditas até o passo atual e a informação aprendida durante o treinamento da rede encarregada pela predição das classes gramaticais. A Figura 4.3 mostra a visão geral do modelo proposto pelo autores.

No trabalho, que não faz uso de qualquer mecanismo de atenção, os autores justificam o uso de POS argumentando que classes gramaticais específicas se correlacionam com informações visuais da imagem de diferentes maneiras, enquanto outras, apesar de não se relacionarem com a imagem a ser descrita, se necessárias para expressas relações semânticas entre palavras a serem preditas.

4.1.2.2 CNN \rightarrow *Attention* \rightarrow RNN

Os autores em (XU et al., 2016) se inspiraram no trabalho da Figura 4.1 para utilizar um mecanismo de atenção com as características extraídas por uma CNN. A arquitetura *Encoder-decoder* é implementada com a diferença de que a representação da imagem

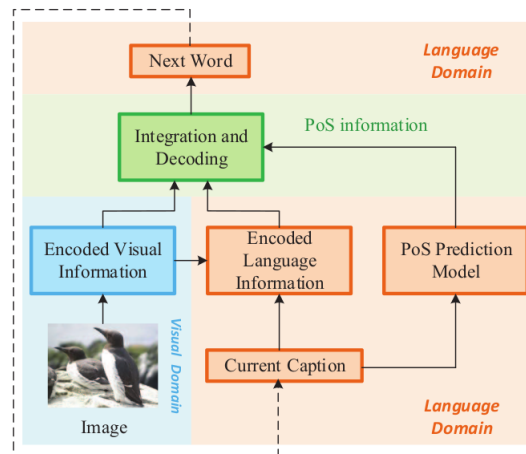


Figura 4.3: Arquitetura do modelo em (ZHANG et al., 2021b). Fonte: (ZHANG et al., 2021b)

feita pelo *Encoder* é adaptada pelos pesos de atenção (calculados com uso da palavra anteriormente gerada) antes de ser utilizada pelo *Decoder* para inferir a saída atual. Utilizando a definição de mecanismo de atenção dada por (VASWANI et al., 2017) para explicar o uso neste trabalho, tem-se:

- requisição (*query*): atual estado interno do *Decoder* sendo comparado com as características da imagem;
- chave (*key*): características (extraídas pela CNN) da imagem de entrada sendo comparadas com a requisição;
- valor (*value*): valor dado para a relevância da característica extraída (chave) e o estado atual (relevância).

Pela Figura 4.4, observa-se que, conforme essa representação da imagem muda, uma parte diferente da imagem é utilizada para inferir cada palavra da sentença final e, de maneira análoga ao que foi visto em (BAHDANAU; CHO; BENGIO, 2014), esse “foco” se relaciona com a intuição humana (a parte usada para predizer a palavra “pássaro”, em azul na Figura 4.4, por exemplo).

Este trabalho foi o primeiro a empregar o mecanismo de atenção para descrição de imagens (BERNARDI et al., 2017). E foi utilizado em nosso primeiro trabalho relacionado na área (GONDIM.; CLARO.; SOUZA., 2022). Uma possível desvantagem é a informação utilizada a cada passo de predição do *Decoder*, palavras como “de”, “em”, “um” (preposições ou determinantes) podem não precisar de informação visual para serem inferidas, mas ainda assim o mecanismo de atenção tenta indicar ao *Decoder* onde, na imagem, deve estar direcionado o foco para gerar palavras que não precisam de tanta informação visual.

O problema de utilizar informações visuais para predizer preposições ou determinantes é abordado no trabalho feito por (LU et al., 2017), no qual os autores argumentam

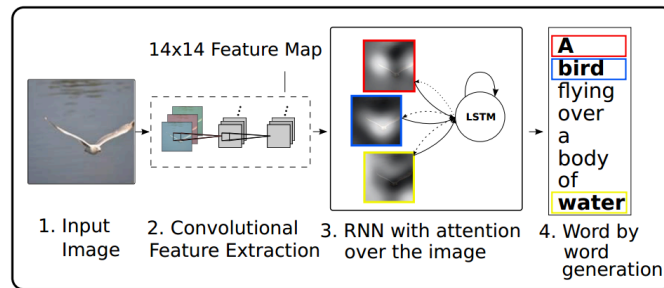


Figura 4.4: Exemplo de geração de descrição em (XU et al., 2016). Fonte: (XU et al., 2016)

que o *Decoder* precisa de pouca ou nenhuma informação visual para prever palavras como “de” ou “o”. Para lidar com esse problema, é proposto o uso de um modelo de atenção adaptativo (*Adaptive Attention Model*). Utilizando uma “sentinela visual”, uma representação da informação armazenada no *Decoder*, a arquitetura proposta é capaz de adaptar a predição da próxima palavra para utilizar apenas informações do modelo de linguagem caso o mecanismo de atenção decida que a informação da imagem não é necessária. A Figura 4.5 exemplifica o uso da sentinela visual, em que o mecanismo de atenção adaptativo decide pelo uso, ou não, de informação visual para prever cada palavra. Palavras em que a informação visual é importante como “homem” (*man*) ou “snowboard” tem alta chance de usar informações da imagem, diferente de “truque” (*trick*) ou “em” (*on*).

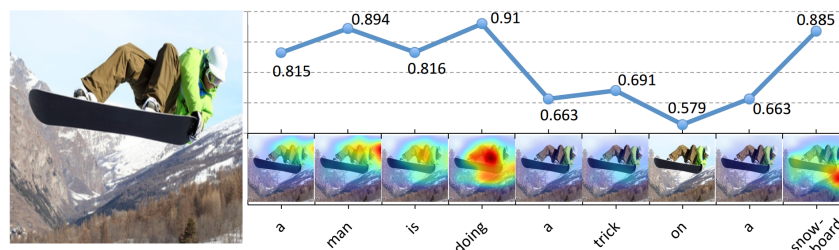


Figura 4.5: Exemplo de geração de descrição em (LU et al., 2017) onde a informação visual é utilizada ou não a depender da palavra predita.

Ainda que seja utilizada a informação de quando usar ou não informações visuais, o modelo empregado ignora a estrutura linguística contida nas sentenças de treinamento. Como o uso de classes gramaticais, que podem auxiliar a definir quando uma informação visual deve ser utilizada, melhorando a capacidade do modelo de “focar”, ou não, na imagem a cada predição de palavra.

4.1.2.3 CNN → Transformer

Com o advento de *Transformers*, a pesquisa em *Image captioning* também passou a

utilizar a arquitetura para criar relações entre texto e imagem que pudessem ser utilizadas para criar descrições. Em (CORNIA et al., 2020), os autores utilizam a arquitetura *Transformer*, modificando a camada de *Self-attention* adicionando um espaço nos vetores de chave e valor, esses espaços são capazes de armazenar informação aprendida ao longo do treinamento, por conta do espaço a mais adicionado, os autores dão o nome de *Memory-Augmented Attention*. *Memory-Augmented Attention* é capaz de aprender informações de relacionamento que não estão presentes nas características extraídas por uma CNN, como por exemplo: tendo duas características, uma de “pessoa” e outra de “bola de basquete”, inferir o conceito de “jogadora” ou “jogo”. Aumentar o espaço dos vetores do mecanismo de *Self-attention* foi a estratégia utilizada pelo autores em (CORNIA et al., 2020) para codificar relações visuais durante o treinamento de modelos, mas essas relações podem ser também aprendidas a partir das classes gramaticais existentes nas sentenças de treinamento.

Em (HUANG et al., 2019), os autores estendem o mecanismo de atenção para determinar se há relevância entre a requisição feita e o resultado do mecanismo de atenção aplicado. Isso se deve ao problema apontado de que, nem sempre a requisição feita ao modelo de atenção tem uma informação útil a ser apontada vinda do *Encoder*, isso leva o *Decoder* a tomar decisões erradas por conta de informações que não são úteis. A solução apresentada é adicionar outro mecanismo de atenção à saída da atenção anterior, obtendo como resultado o conhecimento útil para gerar a descrição. Por conta da montagem da arquitetura, os autores a chamam de *Attention on Attention* (AoA). AoA pode ser aplicada tanto no *Encoder*, para determinar relações entre as características extraídas pelo *Self-attention*, quanto no *Decoder* para filtrar resultados errados do mecanismo de atenção¹

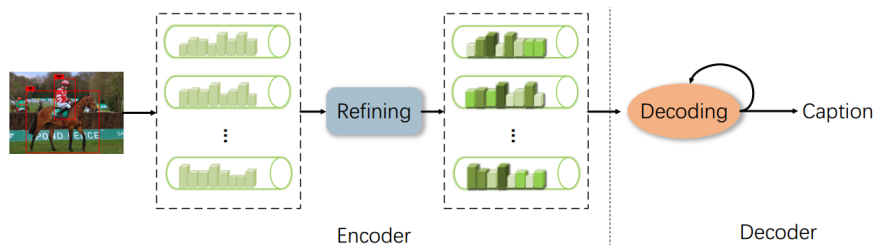


Figura 4.6: Arquitetura do modelo em (HUANG et al., 2019). Fonte: (HUANG et al., 2019)

A Figura 4.6 demonstra a solução do problema apontado no passo de refinamento (*Refining*) entre *Encoder* e *Decoder*, no qual aos vetores de atenção antes do passo de refinamento são dadas as devidas relevâncias para as relações existentes nas características extraídas da imagem de entrada. O mecanismo AoA procura resolver o problema de uma possível falta de relação entre a palavra a ser predita e o contexto retornado pelo modelo de atenção, problema que pode também ser solucionado encontrando relações entre as

¹A arquitetura do trabalho em (HUANG et al., 2019) utiliza *Self-attention* no *Encoder* apenas, contando com RNNs para gerar a sentença da descrição.

classes gramaticais da entrada e as características extraídas da cena.

No trabalho realizado por (PAN et al., 2020), os autores observam que a tarefa de *Image captioning* raramente explora a característica multimodal do problema nos estágios iniciais de treinamento, tratando conteúdos visuais e textuais de maneira separada. Os autores argumentam que o modelo de *Self-attention* explora apenas interações de 1º ordem entre texto e imagem, por conta operação de produto interno entre chave e requisição (no modelo de atenção), e que uma maneira de atenuar esse problema é encontrar interações de maior ordem entre características visuais e textuais. Para isso eles definem um bloco de atenção unificado: *X-Linear attention block*. Utilizando *pooling bilinear*, o *X-Linear attention block* é capaz de extrair interações de 2ª ordem entre a palavra a ser predita e as características da imagem.

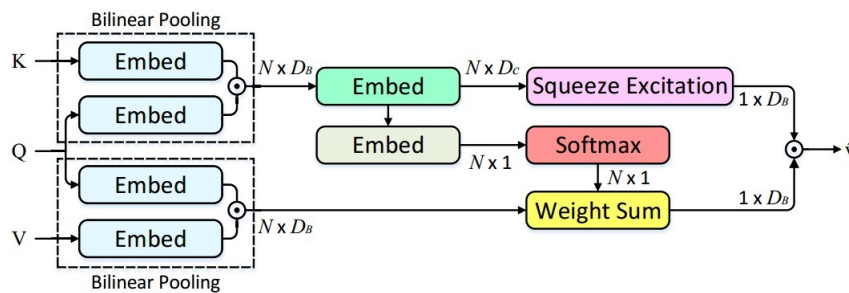


Figura 4.7: Arquitetura do bloco de atenção unificado utilizando *pooling linear*. Fonte: (PAN et al., 2020)

Como ilustrado na Figura 4.7, um bloco de *X-Linear attention* usa *pooling bilinear* nos elementos de entrada do modelo de *Self-attention* para capturar interações de características que sejam mais refinadas, levando a um conjunto melhorado dessas características da imagem de entrada (STEFANINI et al., 2021). O bloco de *X-Linear attention* explora uma limitação do modelo de atenção utilizado em *Self-attention*, fornecendo aos modelos de *Image captioning* uma compreensão mais complexa da relação entre imagem e texto. Ainda assim, não há uso de quaisquer relações ou informações linguísticas que poderiam contribuir com as interações de 2ª ordem sendo criadas.

Mais recentemente, autores em (WANG et al., 2022) afirmam que as predições feitas apenas com informação visual e a sentença parcialmente gerada até então resultam em sentenças genéricas, dada a igualdade no tratamento de palavras visuais e não-visuais durante o treinamento. O artigo propõe um *Transformer* guiado por classes gramaticais, capaz de ajustar o peso dado às informações extraídas da imagem e a informação das palavras sendo preditas. A Figura 4.8 apresenta a arquitetura geral do modelo proposto.

Os autores argumentam que o uso de classes gramaticais como guia na geração de sentenças podem levar a descrições mais refinadas das imagens. Porém o uso de POS é aplicado no topo da arquitetura, após o mecanismo de *Multi-Head Attention* ser aplicado entre informações visuais e o contexto da sentença gerada até o estado atual, o que impossibilita a visualização do “foco” dado pelos mecanismos de atenção e a influência de classes gramaticais no que é retornado.

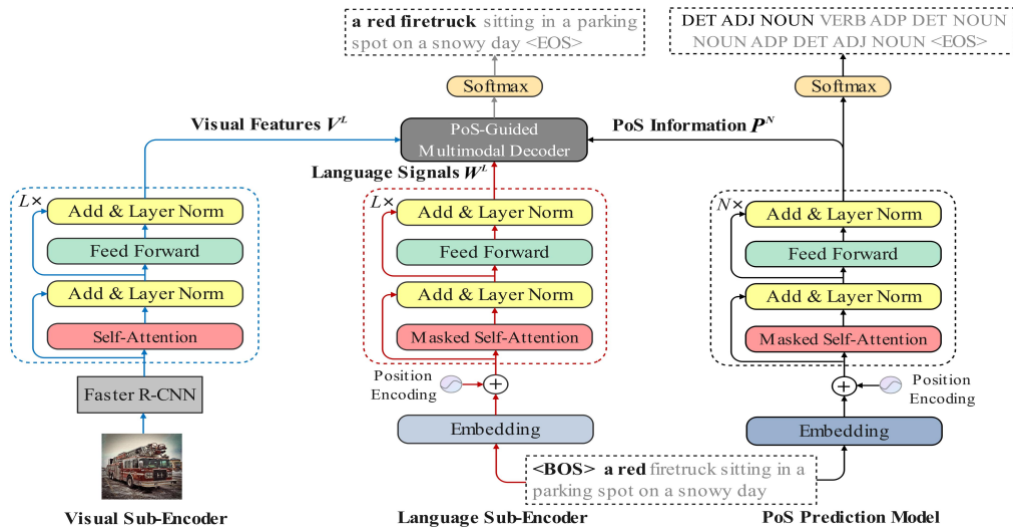


Figura 4.8: Arquitetura do modelo proposto em (WANG et al., 2022). Fonte: (WANG et al., 2022)

4.2 CONJUNTOS DE DADOS

Conjuntos de dados utilizados para *Image captioning* são compostos de imagens e uma lista correspondente de sentenças, que pode ser apenas uma ou múltiplas, embora a existência de múltiplas sentenças para uma imagem ajude na generalização de um modelo, adicionando variabilidade aos dados de treino. As propriedades das sentenças para cada imagem também influenciam nas saídas do modelo a ser treinado, como tamanho da sentença e escolha de palavras para o vocabulário (STEFANINI et al., 2021).

Os conjuntos de dados *Flickr8K* e *Flickr30K* são ambos construídos a partir de imagens encontradas no Flickr², que consistem em cenas comuns ao dia-a-dia. O Flickr8k (HODOSH; YOUNG; HOCKENMAIER, 2013) é um conjunto criado para a tarefa de *Image captioning* que possui 8092 imagens, com cinco descrições cada. Os autores propõem que o conjunto de dados seja utilizado para tratar a tarefa de descrever imagens como um problema de ranqueamento: dada uma imagem, encontrar a sentença que melhor se ajusta a ela. A Figura 4.9 mostra um exemplo de imagem do Flickr8k.

As sentenças de descrição das imagens do conjunto de dados foram anotadas através "crowdsourcing", utilizando serviços de vários anotadores via ferramenta *web* (*Amazon Mechanical Turk*). Aos anotadores, foi solicitado que descrevessem as pessoas, objetos, cenas e atividades mostradas, sem que nenhuma outra informação de contexto fosse adicionada (HODOSH; YOUNG; HOCKENMAIER, 2013). Os autores afirmam que essas instruções resultaram em descrições conceituais, considerando apenas as informações contidas nas imagens unicamente. Os autores do conjunto de dados também fornecem divisões padronizadas para treino, validação e teste.

As sentenças originais anotadas para a Figura 4.9 são:

²<https://www.flickr.com/>



Figura 4.9: Exemplo de imagem do Flickr8k. Fonte: (HODOSH; YOUNG; HOCKEN-MAIER, 2013)

- *A black and white dog is catching a Frisbee in the yard.*
- *A black and white dog is trying to catch a Frisbee in the air.*
- *A dog jumps to catch a red Frisbee in the yard.*
- *Dog is jumping up on a very green lawn to catch a Frisbee.*
- *The black and white dog tries to catch a red Frisbee on green grass.*

Apesar de aparecer pouco em trabalhos mais recentes (que empregam *Transformers*), conjunto de dados *Flickr8k* foi utilizado em dois momentos em suas versões original e traduzida para o português: em experimentos, que serão mostrados no próximo capítulo, e para ajuste de hiper-parâmetros no modelo proposto.

O *Flickr8k* foi estendido no trabalho feito por (YOUNG et al., 2014), no qual os autores apresentam o *Flickr30k*, que possui 31783 imagens com 5 descrições cada, as sentenças que descrevem cada imagem foram criadas seguindo as mesmas instruções do conjunto antecessor. Em seu trabalho, (KARPATHY; FEI-FEI, 2014) dividiram o conjunto de dados em subconjuntos de treino, validação e teste comumente utilizados em outros trabalhos e que também é utilizado neste. No presente trabalho, o *Flickr30k* foi utilizado para análise das saídas dos mecanismos de atenção utilizados em um modelo de *Image Captioning* baseado na arquitetura de *Transformers*.

Atualmente, um conjunto de dados comumente usado para treinamento e avaliação comparativa entre modelos ((HUANG et al., 2019; CORNIA et al., 2020; PAN et al., 2020)) é o *Microsoft COCO Dataset*³ (LIN et al., 2015). O *COCO Dataset* foi desenvolvido para impulsionar o avanço nas pesquisas em três problemas em entendimento de cena (*scene understanding*): detectar visualizações não icônicas (visualizações que remetam a um evento ou localização) de objetos, raciocínio contextual entre objetos e a localização precisa de objetos na cena. Além de anotações de localização de objetos e segmentação⁴, cada uma das 123287 imagens possui 5 sentenças correspondentes.

³<https://cocodataset.org/>

⁴Segmentação da imagem é a tarefa de particionar os pixels da imagem em classes diferentes de objetos (MINAEE et al., 2021).

Uma prática comum para avaliação comparativa entre os trabalhos explorados, é o uso de uma separação padronizada entre treino, validação e teste, definida em (KARPATHY; FEI-FEI, 2014), no qual o conjunto de dados é dividido entre treino (com 82793 imagens), validação (com 5000 imagens), teste (com 5000 imagens) e o restante da validação (com 30504 imagens). Além disso, há um conjunto oficial de teste online com mais de 40 mil imagens em que cada uma possui 40 sentenças que são mantidas privadas⁵.

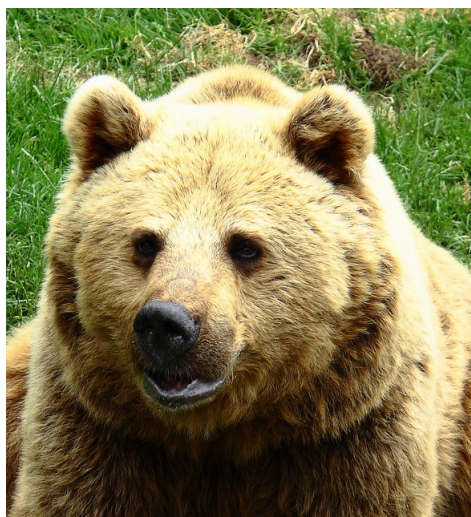


Figura 4.10: Exemplo de imagem do *Microsoft COCO Dataset*. Fonte: (LIN et al., 2015)

Um exemplo de imagem do conjunto de dados é mostrado na Figura 4.10, com as seguintes descrições:

- *A big burly grizzly bear is show with grass in the background.*
- *The large brown bear has a black nose.*
- *Closeup of a brown bear sitting in a grassy area.*
- *A large bear that is sitting on grass.*
- *A close up picture of a brown bear's face.*

Existe pouca pesquisa envolvendo *Image captioning* para o português e isso pode se dar por conta da falta de conjuntos de dados apropriados para a linguagem. O conjunto de dados #PraCegoVer (SANTOS; COLOMBINI; AVILA, 2022) é o primeiro trabalho feito com intenção de juntar imagens e sentenças inteiramente em português. Os autores publicaram dois conjuntos de dados com sentenças em português: #PraCegoVer-63K e #PraCegoVer-173K, com respectivamente 62935 e 173337 imagens. Os conjuntos de dados são desafiadores pois, além de possuírem apenas uma sentença por imagem, suas

⁵<https://competitions.codalab.org/competitions/3221>

anotações não foram feitas para a tarefa de *Image-captioning*, mas para a acessibilidade de pessoas cegas nas redes sociais. Os autores desenvolveram um processo para fazer coleção, pré-processamento e análise de dados para o *Instagram* e buscaram por pares de imagens e sentenças que tivessem a marcação “#PraCegoVer”⁶. As sentenças diferem bastante (em tamanho) dos conjuntos anteriores (Flickr8k, Flickr30k e *Microsoft COCO Dataset*). Os autores fazem uma comparação adaptada na Tabela 4.1, mostrando o tamanho dos conjuntos de dados, as quantidades de palavras únicas, o tamanho médio das sentenças e o desvio padrão do tamanho das sentenças.

Conjunto de Dados	Tamanho	Quantidade de palavras únicas	Tamanho médio sentença	Desvio padrão
Flickr8k	8091	8488	10,8	3,8
Flickr30k	31014	17984	12,3	5,2
<i>Microsoft COCO Dataset</i>	123287	24940	10,6	2,2
#PraCegoVer-63K	62935	55029	37,8	26,8
#PraCegoVer-173K	173337	93085	39,3	29,7

Tabela 4.1: Tabela com os dados dos *datasets* COCO, Flickr8k, #PraCegoVer-63K e #PraCegoVer-173K. Fonte: adaptado de (SANTOS; COLOMBINI; AVILA, 2022).

A Figura 4.11 ilustra um exemplo de imagem presente no #PraCegoVer, a sentença desta imagem é: “Em um ambiente externo, uma gata de pelagem branca e caramelo, está deitada de olhos fechados, sua ninhada de filhotinhos esta ao seu redor, um dos gatinhos esta olhando para a câmera, ele tem olhos cor de mel e pelagem branca, preta e caramelo, o restante dos filhotes, estão desfocados. No canto inferior direito, está escrito, “PremieRpet” em letras alaranjadas.”



Figura 4.11: Exemplo de imagem do #PraCegoVer. Fonte: (SANTOS; COLOMBINI; AVILA, 2022)

Os conjuntos de dados #PraCegoVer são desafiadores e possuem sentenças naturalmente coletadas em português. Além disso, os autores treinaram o modelo (HUANG et al., 2019) com as imagens coletadas, o que fornece uma avaliação comparativa na língua portuguesa.

⁶“#PraCegoVer” é um movimento que procura melhorar a acessibilidade de pessoas com deficiência visual ao acesso da internet incluindo descrição de imagens em postagens em redes sociais e sites (SANTOS; COLOMBINI; AVILA, 2022).

4.3 RESUMO DO CAPÍTULO

Neste capítulo abordou-se trabalhos relacionados da área de *Image Captioning*. Primeiro apresentando os avanços na área, desde as abordagens iniciais, com o uso de regras e *templates*, até a implementação de redes neurais profundas com arquiteturas do tipo “Encoder-decoder”, que, nos trabalhos estado da arte, são implementadas por meio de *Transformers*. Em seguida, são abordados conjuntos de dados comumente utilizados para treinamento e comparação de modelos, onde observa-se a dominância de conjuntos feitos para a língua inglesa, havendo apenas um trabalho conhecido com sentenças na língua portuguesa.

LINGUISTIC IMAGE CAPTIONING

O trabalho proposto tem como objetivo principal “*Desenvolver uma arquitetura para descrever imagens em Português utilizando informações linguísticas que sejam capazes de melhorar a geração de textos*”. O uso de informações linguísticas será explorado adicionando as classes gramaticais (POS) das sentenças no fluxo de treinamento. Além disso, nosso trabalho pretende focar na língua portuguesa durante o desenvolvimento, utilizando conjuntos de dados apropriados para a tarefa.

5.1 METODOLOGIA

Esta seção introduz os conjuntos de dados utilizados, a extração de classes gramaticais utilizadas e os métodos de avaliação empregados.

5.1.1 Conjuntos de Dados

O primeiro Objetivo Específico do trabalho (OE 1) é analisar um conjunto de dados traduzido para dar seguimento aos experimentos que serão realizados, visto que não há um conjunto de dados para português. As traduções foram realizadas empregando a ferramenta LibreTranslate¹, que faz tradução automática de sentenças utilizando OpenNMT (KLEIN et al., 2017), a ferramenta foi utilizada localmente para tradução. A escolha de palavras e os tamanhos das sentenças influenciam no modelo final treinado, portanto, para cada conjunto de dados traduzido, a quantidade de palavras únicas (vocabulário), o tamanho médio de sentenças e o desvio padrão desse tamanho médio são analisados.

O primeiro conjunto de dados traduzido foi o *Flickr8k* (HODOSH; YOUNG; HOCKENMAIER, 2013). Cada imagem contém cinco sentenças e cada uma delas foi traduzida utilizando o OpenNMT (KLEIN et al., 2017). Um exemplo são as sentenças da Figura 5.1 traduzidas:

- Um cão preto e branco está pegando um Frisbee no pátio.

¹<https://github.com/LibreTranslate/LibreTranslate>



Figura 5.1: Exemplo de imagem do Flickr8k. Fonte: (HODOSH; YOUNG; HOCKENMAIER, 2013)

- Um cão preto e branco está tentando pegar um Frisbee no ar.
- Um cão salta para pegar um Frisbee vermelho no pátio.
- O cão está pulando em um gramado muito verde para pegar um Frisbee.
- O cão preto e branco tenta pegar um Frisbee vermelho na grama verde.

Após a tradução, o vocabulário que possuía 8488 palavras aumentou para 9780, o tamanho médio de sentenças foi alterado de 10,8 para 11,16 e o desvio padrão de 3,8 para 4,00. Além disso, alguns erros de tradução surgiram, como por exemplo a sentença “Um cão preto está saltando sobre um *log* ao longo de uma praia” em que a palavra *log* não é traduzida. Além disso, na sentença “*A dirt bike racer jumps over a slope.*” a tradução é “Um motociclista de sujeira salta sobre uma inclinação”. Esses erros de tradução são inseridos durante o treinamento dos modelos nos experimentos realizados, representando uma limitação na criação de modelos para a língua portuguesa.

O segundo conjunto de dados utilizado é o *Flickr30k* (YOUNG et al., 2014). A Figura 5.2 mostra um exemplo de uma das imagens contidas.



Figura 5.2: Exemplo de imagem do *Flickr30k*. Fonte: (YOUNG et al., 2014)

A tradução automatizada das sentenças de referência da Figura 5.2 é mostrada abaixo:

- Um homem está flexionando seus biceps enquanto está em um telhado com várias chaminés em segundo plano;
 - A man is flexing his biceps while standing on a rooftop with multiple chimneys in the background
- Um homem em uma camisa marrom e listrada está flexionando seus braços em cima de uma estrutura em um telhado;
 - A man in a brown, striped shirt stands flexing his arms on top of a structure on a rooftop
- Um homem de camisa listrada e jeans sujos está triunfantemente em cima de um telhado;
 - A man in a striped shirt and dirty jeans is triumphantly standing on top of a roof
- Um homem burly em cima de tubos de plástico flexionando em uma camisa polo preto e branco;
 - A burly man on top of plastic pipes flexing in a black and white polo shirt
- Um tipo a flexionar os músculos.
 - One guy flexing his muscles

Assim como no caso visto com o conjunto de dados *Flickr8k*, a tradução automatizada das sentenças acarretou em mudanças no tamanho do vocabulário e no tamanho médio das sentenças. A Tabela 5.1 demonstra as características dos conjuntos de dados utilizados e onde cada um deles foi aplicado. Para facilitar a separação dos conjuntos, adicionou-se aos seus nomes os finais “-EN” (para indicar o conjunto de dados original na língua inglesa) e “-PT” (para indicar o conjunto de dados com sentenças automaticamente traduzidas para o português).

Conjunto de Dados	Tamanho	Quantidade de palavras únicas	Tamanho médio sentença	Aplicação
Flickr8k-EN	8091	8488	10,8 ± 3,8	Experimentos 1 e 2
Flickr8k-PT	8091	9780	11,2 ± 4,0	Experimentos 1, 2 e 4
Flickr30k-EN	31014	17984	12,3 ± 5,2	Experimento 3
Flickr30k-PT	31014	21367	12,7 ± 5,5	Experimentos 3 e 4

Tabela 5.1: Tabela com informações dos conjuntos de dados utilizados nos experimentos.

5.1.2 Extração de Classes gramaticais

A arquitetura a ser proposta na próxima seção faz uso de classes gramaticais para guiar o treinamento do modelo. Essas informações linguísticas foram previamente extraídas com a biblioteca *spaCy* (MONTANI et al., 2022), através do seu modelo *large*.

Após extrair as classes gramaticais das sentenças da Figura 5.2, tem-se as seguintes anotações:

- DET NOUN AUX VERB DET NOUN ADV AUX ADP DET NOUN ADP DET NOUN ADP ADJ NOUN;
 - Um homem está flexionando seus biceps enquanto está em um telhado com várias chaminés em segundo plano
- DET NOUN ADP DET NOUN ADJ CCONJ VERB AUX VERB DET NOUN ADP ADP ADP DET NOUN ADP DET NOUN;
 - Um homem em uma camisa marrom e listrada está flexionando seus braços em cima de uma estrutura em um telhado
- DET NOUN ADP NOUN VERB CCONJ NOUN ADJ AUX ADV ADP ADP ADP DET NOUN;
 - Um homem de camisa listrada e jeans sujos está triunfantemente em cima de um telhado
- DET NOUN NOUN ADP ADP ADP NOUN ADP NOUN VERB ADP DET NOUN NOUN ADJ CCONJ ADJ;
 - Um homem burly em cima de tubos de plástico flexionando em uma camisa polo preto e branco
- DET NOUN SCONJ VERB DET NOUN.
 - Um tipo a flexionar os músculos

Os erros de tradução, se refletem na extração de classes gramaticais: na quarta sentença, a palavra “burly” (que não foi traduzida) é erroneamente classificada como NOUN (“substantivo” como visto na Tabela 3.2), uma vez que sua tradução, “corpulento” ou “forte”, se encaixa como ADJ (Adjetivo).

As classes gramaticais extraídas foram organizadas de maneira similar às sentenças, sendo cada conjunto com 5 anotações atribuído à sua respectiva imagem. Durante a fase de treinamento da arquitetura proposta na seção 5.2, as classes gramaticais serão passadas juntamente com as respectivas sentenças e imagens.

5.1.3 Avaliação

As avaliações das sentenças geradas pelos modelos treinados foram feitas de forma qualitativa e quantitativa. Primeiramente, foram analisadas as sentenças geradas por um modelo de *Image captioning* utilizando uma CNN, um mecanismo de atenção e uma RNN (similar ao trabalho feito em (XU et al., 2016)). Essa análise serviu de base para seleção de erros comuns utilizados em avaliação qualitativa. Já a avaliação quantitativa, a ser realizada pra avaliar o desempenho da arquitetura proposta, foi feita com métricas automatizadas comuns às utilizadas nestes trabalhos: BLEU 1-4 (PAPINENI et al., 2002) e METEOR (BANERJEE; LAVIE, 2005)².

²As métricas foram calculadas utilizando a ferramenta <https://github.com/tylin/coco-caption>

5.2 ARQUITETURA PROPOSTA

A arquitetura utiliza informações extraídas previamente das sentenças do corpus utilizado: classes gramaticais das palavras, por conta disso, foi dado o nome de *Linguistic Image Captioning (LIC)*. O uso de classes gramaticais durante o treinamento do modelo segue conclusões dos trabalhos citados no capítulo 4. Informações como a falta de necessidade de informações visuais para prever palavras não-visuais (LU et al., 2017), ou o uso de POS como guia durante o treinamento de RNNs (HE et al., 2017) e que palavras de diferentes classes gramaticais se correlacionam de maneiras diferentes com o que é apresentado na imagem (ZHANG et al., 2021b), levaram à utilização de informações sintáticas para tentar melhorar a interação entre imagem e texto dentro do mecanismo de *Multi-Head Attention*.

A arquitetura é demonstrada na Figura 5.3, e tem inspiração do modelo neural de um *Transformer* com o uso de classes gramaticais para guiar o treinamento do mecanismo de *Self-attention*.

Na arquitetura *Encoder-decoder* com *Transformer* e CNN, o uso de *Self-attention* é aplicado nas características extraídas da imagem, extraindo relações entre regiões desta entrada que são utilizadas no *Decoder* para predição de cada palavra. O LIC adiciona um etapa a mais após a codificação da imagem, as informações retornadas pelo *Encoder* são utilizadas para computar a relevância entre as informações da imagem e as classes gramaticais das palavras até a posição atual sendo predita no bloco chamado de *Decoder POS*. Com esse passo a mais, pretende-se guiar o modelo final especificando uma parte da imagem dependendo da classe gramatical da palavra anterior. De maneira análoga ao procedimento de predição, sem usar classes gramaticais, esse passo do LIC é semelhante a “prever uma sequência de classes gramaticais dada uma imagem”. Durante a etapa de predição, o *Decoder POS* é responsável por inferir a sequência de classes gramaticais para uma dada imagem, que é utilizada pelo *Decoder* para geração da sentença final. É importante que o *Decoder POS* seja treinado concomitantemente com o modelo pois as informações das classes gramaticais, que foram extraídas com uma ferramenta externa, não estarão disponíveis durante o passo de predição.

O *Decoder* é implementado adicionando mais duas subcamadas de *Masked Multi-Head Attention* com as operações de conexão residual e normalização (destacadas com linhas pontilhadas na Figura). A primeira recebe as classes gramaticais da sequência de entrada (inicialmente extraídas com o *spaCy* e, após o treinamento, informadas pelo *Decoder POS*) e gera uma representação desse vetor, utilizando o mesmo como requisição, chave e valor, portanto realizando uma operação de *Self-attention*. Na segunda subcamada, a representação gerada é utilizada como chave a ser comparada com a saída dada pelo primeiro bloco de *Masked Multi-Head Attention*, que gerou uma representação da sentença original, e é utilizada como requisição e valor, portanto retornando a relevância entre a representação das classes gramaticais e a representação da sentença dada como entrada. Por fim, esse vetor é passado para a subcamada de *Multi-Head Attention* que irá informar onde o modelo deve “focar” na imagem para fazer a predição da próxima palavra. Para treinar o modelo neural, foram utilizados os conjuntos de dados *Flickr8k* para seleção de hiper-parâmetros e *Flickr30k* para treinamento final, juntamente com as

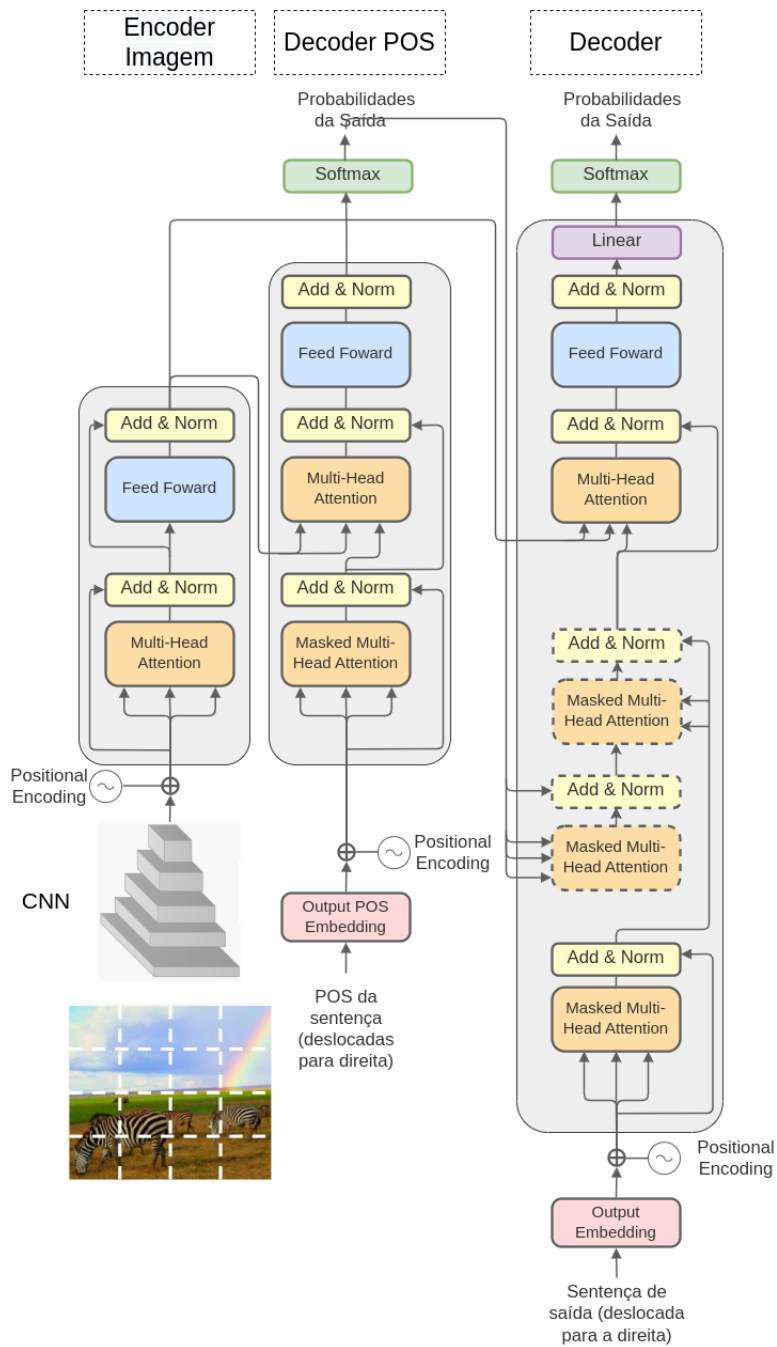


Figura 5.3: Arquitetura do modelo de descrição linguística de imagens do trabalho Fonte: adaptado de (VASWANI et al., 2017).

classes gramaticais extraídas de cada um deles.

5.3 EXPERIMENTOS

Três experimentos foram realizados com as saídas obtidas de diferentes modelos de *Image captioning* treinados. Primeiro avaliou-se o desempenho alcançado por uma mesma arquitetura com conjuntos de dados diferentes: um com as sentenças traduzidas do inglês para o português e outro com as sentenças originais em inglês. Em seguida, erros encontrados nas sentenças geradas foram categorizados para serem identificados por anotadores externos. No terceiro experimento, o conjunto de dados *Flickr30k* foi utilizado para treinar um modelo de *Image captioning* utilizando um *Transformer* (como mostrado na Figura 3.14) e avaliar o comportamento das saídas dos mecanismos de atenção de acordo com mudanças morfológicas.

Nos experimentos 1 e 2, para selecionar os dois modelos utilizados, primeiro avaliou-se quantitativamente quatro arquiteturas de *Image captioning* treinadas com o conjuntos de dados *Flickr8k* com sentenças traduzidas para o português. Em seguida uma das arquiteturas foi selecionada e treinada com as mesmas configurações, porém utilizando o *Flickr8k* original (com sentenças em inglês). Após isso, as métricas das duas arquiteturas similares, cada uma treinada com uma linguagem diferente, foram comparadas para todas as sentenças geradas em cada uma das mil imagens do conjunto de dados de teste.

Os treinamentos foram realizados utilizando um modelo adaptado de (XU et al., 2016), ilustrado na Figura 5.4. O modelo utiliza arquitetura *Encoder-decoder*, com uso de um mecanismo de atenção definido em (BAHDANAU; CHO; BENGIO, 2014) e um no modelo de extração de características mais recente: EfficientNetB7 (TAN; LE, 2020). Primeiro extraímos características e então, ao gerar cada palavra, a RNN utiliza o mecanismo de atenção para “focar” em uma parte das características extraídas da imagem.

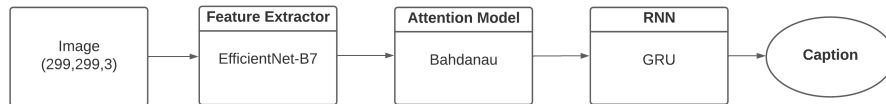


Figura 5.4: Arquitetura simplificada do modelo de descrição de imagens utilizado em (GONDIM.; CLARO.; SOUZA., 2022) Fonte: preparado pelos autores.

5.3.1 Experimento 1

A primeira análise quantitativa feita, foi por meio de comparação entre métricas BLEU 1 a 4 e METEOR. Os quatro modelos treinados, com configurações similares de extração de características, mecanismo de atenção e RNN, se diferenciavam apenas no tamanho do espaço de representação de palavras (*Embedding*). Dois modelos tinham um *Embedding* com espaço de dimensões de tamanho 300, um deles inicializado com pesos pré-treinados com GloVe (PENNINGTON; SOCHER; MANNING, 2014), enquanto que o outro foi inicializado com valores aleatórios e treinado juntamente com os dados das sentenças de treinamento. Em outros dois modelos, foram empregados 600 dimensões para o *Embedding*, sendo um deles inicializado com GloVe e o outro treinado similarmente ao de

tamanho 300. Os treinamentos foram realizados por 50 épocas, com taxa de aprendizagem de 0,001 e com lotes de 128 imagens. Cada imagem foi redimensionada para 299 *pixels* de altura por 299 *pixels* de largura.

5.3.2 Resultados 1

Os resultados são apresentados na Tabela 5.2

Embedding	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR
300d Embedding	41,24	24,25	13,77	7,89	14,81
600d Embedding	41,18	24,53	14,29	8,33	15,3
GloVe 300	39,96	23,74	13,97	8,23	15,46
Glove 600	41,11	24,77	14,32	8,24	14,81

Tabela 5.2: Métricas obtidas com o modelo treinado com o conjunto de dados em Português.

Os resultados apresentados mostraram métricas não muito divergentes em relação às dimensões. Portanto o modelo treinado com GloVe e *Embedding* de 300 dimensões, o segundo melhor na comparação entre métricas, foi selecionado para dar continuidade ao experimento. A escolha foi feita por conta do uso de pesos pré-treinados do GloVe que podem levar a um risco menor de *overfitting*³ do modelo. Utilizamos as mesmas configurações do modelo escolhido para treinamento com o conjunto de dados original do Flickr8k (sentenças em inglês). As métricas obtidas são apresentadas na Tabela 5.3.

Embedding	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR
GloVe 300	42,88	25,65	14,66	8,15	15,19

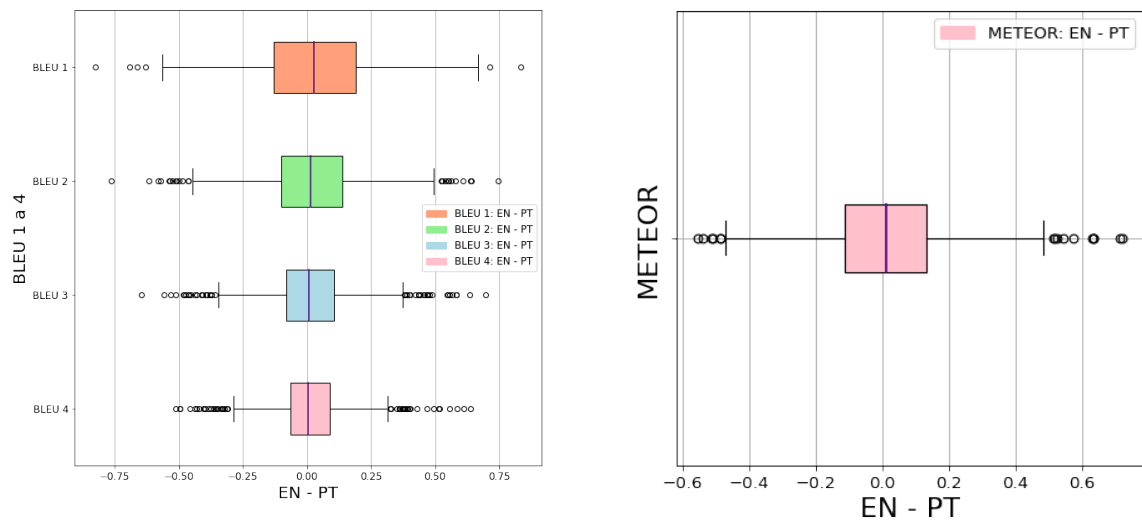
Tabela 5.3: Métricas obtidas do modelo treinado com o conjunto de dados em Inglês.

Os dois modelos finais, com as mesmas configurações e treinados com conjuntos de dados diferentes (um em português e outro em inglês), foram utilizados para comparar as métricas das 1000 imagens de teste do Flickr8k. Para cada uma dessas 1000 imagens, calculou-se as métricas BLEU de 1 a 4 e METEOR para as sentenças em inglês e português e depois foi calculada a diferença entre as métricas para cada imagem. A Figura 5.5 mostra os *boxplots* comparativos dessas diferenças para cada uma das métricas.

Na Figura 5.5a, observa-se que os quantis dos *boxplots* se aproximam com o aumento dos *n-grams* (de 1 a 4), assim como a quantidade de *outliers* também aumenta. Isso pode ser um indicativo de que as métricas tornam-se menos dispersas (mais similares) para os modelos treinados em inglês e em português, conforme o cálculo BLEU fica mais rigoroso. Porém também observa-se (pela crescente quantidade de *outliers*) que há muito

³*Overfitting* indica que o modelo aprendeu representações específicas para os dados de treinamento, performando mal nos dados de validação ou teste (CHOLLET, 2017)

mais valores discrepantes conforme a quantidade de n -grams utilizadas aumenta. A Figura 5.5b mostra as diferenças com a métrica METEOR que possui primeiro e terceiro quantis não muito separados, mas com mínimo e máximo (*whiskers*) mais afastados. Uma observação importante é que, em todos os *boxplots*, houve uma leve assimetria positiva, indicando métricas maiores para os modelos treinados em inglês, isso pode ser devido ao conjunto de dados traduzido que traz alguns erros inerentes do processo de tradução.



(a) Comparação de métricas BLEU 1 a 4

(b) Comparação de métricas METEOR

Figura 5.5: Comparações entre métricas de descrições feitas com as mesmas imagens e linguagens distintas.

5.3.3 Experimento 2

Após essa comparação de métricas, uma avaliação qualitativa foi realizada, através de um formulário enviado para anotadores humanos com 100 imagens aleatoriamente selecionadas. Nesse formulário, os anotadores marcaram a presença dos seguintes erros:

- **Erro 1:** Erro na descrição do personagem (gênero, quantidade, faixa etária, ação etc);
- **Erro 2:** Erro na cor de um objeto;
- **Erro 3:** Erro no cenário (praia, lagoa, montanha);
- **Erro 4:** Erro nos objetos da cena;
- **Erro 5:** Frase mal estruturada (verbo não concorda com sujeito, repetição de palavras);

Além disso, foi perguntado se:

- **Correta:** A sentença gerada descreve a imagem corretamente;
- **Incorreta:** A sentença gerada está errada e não reflete a imagem apresentada.

5.3.4 Resultados 2

As respostas foram contabilizadas e os resultados são apresentados na Figura 5.6. O erro mais presente nas marcações dos anotadores está relacionado à descrição dos personagens da cena, seguido da sinalização de uma frase mal estruturada. Observa-se que os erros mais presentes (erro na descrição do personagem e frase mal estruturada) se relacionam com a QP 4 (“Auxiliar o treinamento de um sistema de Image Captioning com informação das classes gramaticais de cada palavra melhora a sentença gerada?”) uma vez que as informações das classes gramaticais das palavras a serem preditas podem dar suporte para uma descrição menos genérica, bem como auxiliar provendo uma estrutura sintática para a sentença sendo inferida.

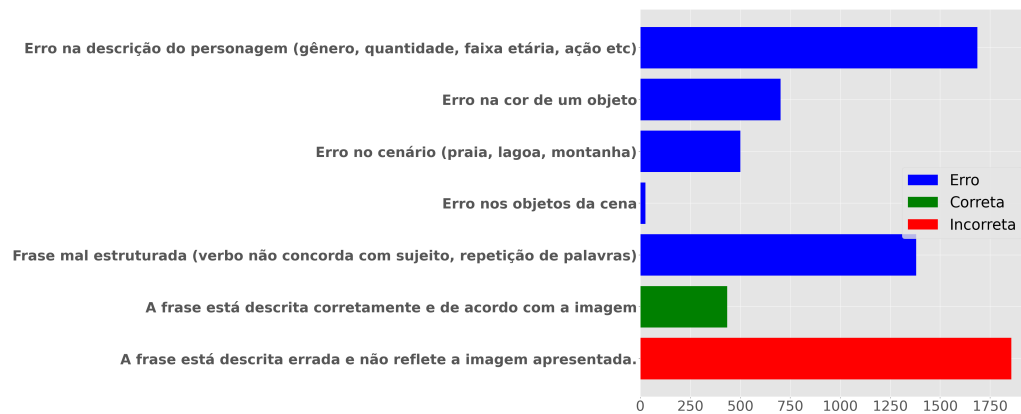


Figura 5.6: Contagens de marcações feitas nos formulários. Fonte: preparado pelos autores.

A pontuação de concordância dos anotadores para cada erro também foi calculada, os valores são apresentados na Tabela 5.4.

A tabela mostra que, embora os Erros 1 e 5 (descrição de personagens e sentença mal estruturada respectivamente) sejam os mais votados, a concordância entre os anotadores é baixa, indicando que pode não ser simples para os anotadores diferenciar o que caracteriza um Erro 1 ou 5. Ao contrário do que é visto para o Erro 4 que tem uma concordância alta quando marcado, apesar de ser o menos sinalizado dentre as imagens, sugerindo que não é constante a presença de objetos errados na descrição de uma cena.

Ainda assim o erro na descrição de cor de objetos (Erro 2) tem uma concordância maior e maior presença de acordo com a Figura 5.6, o que sugere um problema na descrição de alguns objetos. Este problema pode ser explorado com o uso de classes gramaticais para guiar o treinamento de um modelo, por exemplo, o modelo pode utilizar

Erro	Pontuação
Erro 1: Errou descrição do personagem (gênero, quantidade, faixa etária, ação etc)	21,01%
Erro 2: Erro na cor de um objeto	42,96%
Erro 3: Errou o cenário (praia, lagoa, montanha)	23,31%
Erro 4: Errou os objetos da cena	75,58%
Erro 5: Frase mal estruturada (verbo não concorda com sujeito, repetição de palavras)	27,81%
Correta: A sentença gerada descreve a imagem corretamente	36,74%
Incorreta: A sentença gerada está errada e não reflete a imagem apresentada	19,52%
Concordância geral	40,52%

Tabela 5.4: Pontuações de concordâncias entre anotadores para cada erro.

a informação da predição da classe gramatical “substantivo” anteriormente predita para guiar a descrição de um objeto predito anteriormente, focando em informações extraídas que sejam relevantes para o objeto.

5.3.5 Experimento 3

O experimento 3 teve como objetivo avaliar as saídas dos mecanismos de atenção de um modelo de *Image Captioning*, avaliadas de acordo com mudanças morfológicas (de gênero e de modo verbal). Ao analisar o bloco de *Multi-Head Attention* do *Decoder*, é possível visualizar em que parte da imagem o modelo “focou” ao fazer a predição de cada palavra. No caso de um *Transformer*, sua arquitetura, com múltiplas *heads* de atenção, possibilita diferentes análises, uma para cada *heads*.

Autores em (CLARK et al., 2019) e (VIG, 2019) argumentam que ao analisar as saídas de mecanismos de atenção é possível identificar padrões de comportamento. Esses padrões permitiram que os autores: encontrassem vieses contidos nos conjuntos de dados utilizados, localizassem *heads* relevante para estudar um fenômeno específico e identificassem que *heads* diferentes se especializam em certas noções linguísticas durante o treinamento. Porém esses estudos foram realizados com *Transformers* sendo utilizados em tarefas “texto-para-texto”, ou seja, tarefas com entradas e saídas textuais.

O modelo usado para este experimento foi um *Transformer* com duas camadas para o *Encoder* e o *Decoder*, cada mecanismo de *Multi-Head Attention* com 6 *heads*, o *Embedding* foi configurado com tamanho 512 e os módulos de rede neural com conexão para frente (*Feed Forward*) com tamanho de 2048. Os hiper-parâmetros de treinamento utilizados foram lotes de 64 imagens (redimensionadas para o tamanho de 384 *pixels* de altura por 384 *pixels* de largura); tamanho fixo de sentenças de 25 palavras; taxa de aprendizagem com valor de 10^{-5} após 20000 passos de aquecimento. O treinamento ocorreu até que não houvesse melhora na função de perda por 5 épocas seguidas.

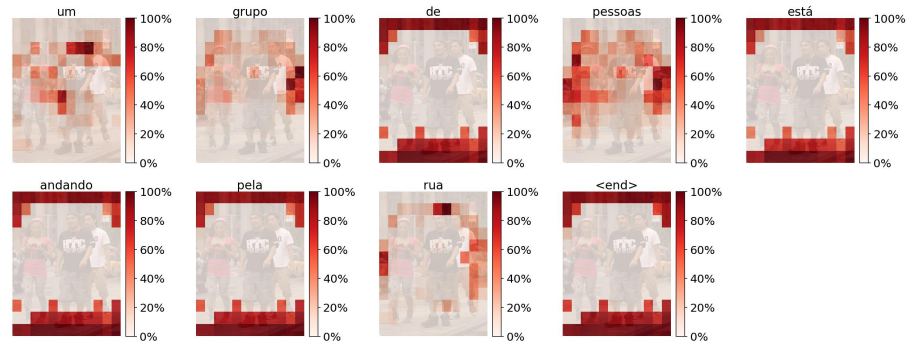
Nenhuma técnica para melhorar as sentenças geradas foi utilizada, como o uso de *beam search*, onde, concomitantemente, mais de uma sentença é gerada para uma imagem sendo escolhida aquela com maior probabilidade total dentre as palavras preditas, ou *Self-critical Sequence Training* onde o modelo, após o treinamento, é otimizado para a métrica CIDEr. O intuito dessas decisões é o de analisar um modelo de *Image Captioning* “puro”, de modo em que apenas as saídas de atenção influenciem na predição final.

Assim como o experimento 1, o mesmo modelo foi utilizado para duas versões do conjunto de dados *Flickr30k*: Flickr30k-EN e Flickr30k-PT; que, respectivamente, deram origem aos modelos Modelo-EN e Modelo-PT.

Durante o processo de geração de sentenças, os pesos de cada mecanismo de atenção (que são interpretados como o “foco” dado em partes da imagem) foram coletados a cada palavra sendo predita. A Figura 5.7 demonstra o método, onde a sentença “um grupo de pessoas está andando pela rua” foi gerada para descrever a imagem mostrada na Figura 5.7a. A Figura 5.7b mostra os pesos de atenção da *head* 1 para cada palavra predita, que, após serem retornados pelo mecanismo de *Multi-Head Attention*, foram distribuídos em uma grade de tamanho 12x12 (este tamanho é definido pela saída da CNN utilizada como extrator de características e pelo tamanho da imagem passada para o modelo) onde cada bloco recebe o correspondente valor do “foco” dado. As saídas do mecanismo de atenção possuem valores diferentes, o somatório desses valores é 1. Para facilitar a visualização, os valores foram normalizados com o maior valor de um bloco correspondendo a 100% na escala ao lado de cada imagem.



(a) Imagem original.
Fonte: (YOUNG et al., 2014).



(b) Sentença gerada e o “foco” dado a cada parte da imagem no momento em que cada palavra foi predita.

Figura 5.7: Exemplo de predição de sentença e do “foco” dado pela *head* 1.

Cada peso retornado é agrupado palavra a palavra e por cada uma das 6 *heads*. Em seguida, esses grupos são sumarizados pela média, possibilitando um investigação da média do foco dado por cada *head*. Um exemplo é mostrado na Figura 5.8, onde a média de atenção dada por cada *head*, quando Modelo-PT prediz a palavra “andando”, é exibida.

Como visto na Figura 5.8, *heads* diferentes apresentam comportamentos médios diferentes ao predizer a palavra “andando”, porém duas delas têm médias diferentes das demais. Enquanto as *heads* 2,3,4 e 6 “atentam” de maneira centralizada na imagem, as demais distribuem os pesos nos cantos das imagens. Nota-se, pela Figura 5.9, que o mesmo não acontece ao realizar o procedimento com a palavra “walking”, com pesos retornados pelo Modelo-EN.

A Figura 5.10 mostra que, ao analisar o mesmo verbo (“andando”), porém em sua forma infinitiva, “andar”, é identificado que o comportamento de 5 *heads* é similar ao

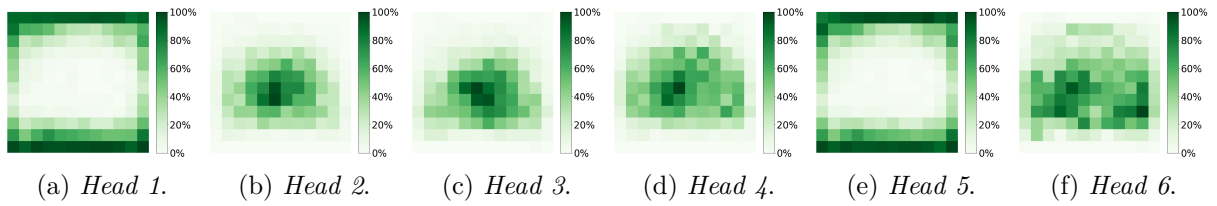


Figura 5.8: Médias de cada *head* de atenção quando a palavra “andando” foi predita pelo Modelo-PT.

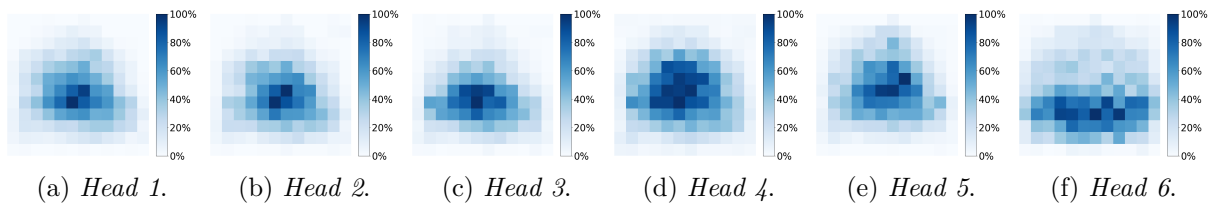


Figura 5.9: Médias de cada *head* de atenção quando a palavra “walking” foi predita pelo Modelo-EN.

visto anteriormente, porém uma delas se altera no funcionamento médio.

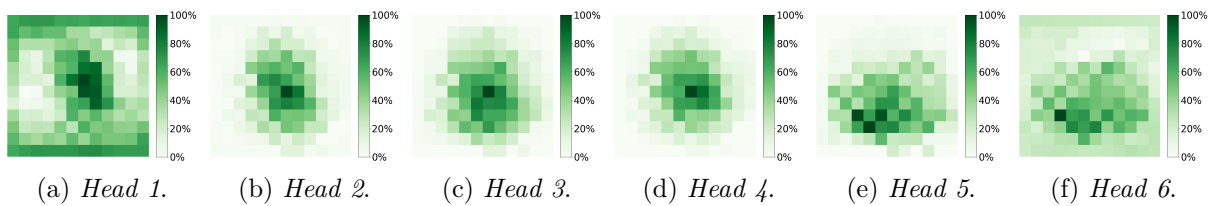


Figura 5.10: Médias de cada *head* de atenção quando a palavra “andar” foi predita pelo Modelo-PT.

Para facilitar a visualização de padrões de comportamento, as diferenças entre as grades, que representam a parte da imagem com maior atenção, são exibidas. Primeiro, cada bloco com o valor médio de atenção para uma palavra foi normalizado pelo valor máximo dos blocos de cada *head*, de forma que o valor máximo de um bloco seja 1, em seguida, calculou-se o módulo das diferenças entre cada bloco normalizado com o respectivo bloco da outra palavra a ser comparada. Desta forma, quando o “foco” retornado pelo mecanismo de atenção em um bloco da grade coincide entre palavras diferentes, o valor da diferença será 0.

A Figura 5.11 exemplifica com a diferença das distribuições de atenção entre as palavras “andando” e “andar” para cada *head*, ordenadas da 1 até a 6. A média e o desvio padrão dessas diferenças são apresentados acima de cada *head* para evidenciar diferenças altas ou baixas.

As diferenças apresentadas são mais perceptíveis justamente na *head* 5, onde foi identificado anteriormente um contraste maior no padrão de comportamento do “foco” da

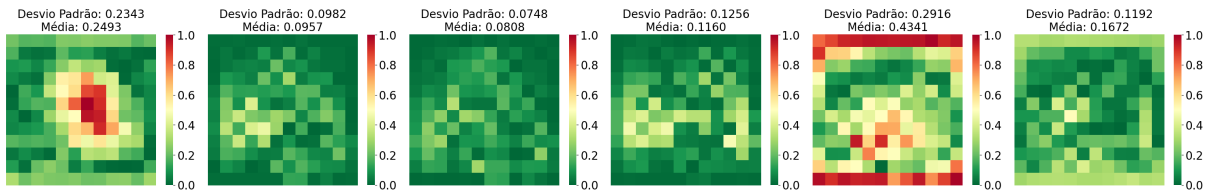


Figura 5.11: Diferenças entre as *heads* do Modelo-PT ao prever as palavras “andando” e “andar”.

atenção. Os resultados desse experimento são exibidos com base nessas mudanças de comportamento.

5.3.6 Resultados 3

Após o treinamento, as 1000 imagens do subconjunto de testes tiveram suas descrições geradas automaticamente. A Tabela 5.5 apresenta os valores obtidos nas métricas BLEU 1-4 e METEOR.

Model	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR
Modelo-EN	60.96	43.35	30.07	20.89	19.56
Modelo-PT	57.80	39.52	26.70	17.97	18.16

Tabela 5.5: Métricas obtidas para cada modelo treinado e avaliado com seus respectivos conjuntos de dados.

A tabela mostra que, apesar de erros de tradução presentes no conjunto de dados Flickr30k-PT, os resultados não diferem muito entre os modelos, permanecendo na casa das unidades para cada métrica avaliada. Esses valores corroboram com o que foi apresentado nos resultados do experimento 1, em que as métricas dos modelos treinados com sentenças originais apresentam valores maiores, mas com diferenças geralmente pequenas quando comparadas.

As diferenças mostradas da média dos pesos de atenção quando as palavras “andando” e “andar” levam aos primeiros resultados, onde é avaliada a flexão verbal através da substituição do sufixo “-r” (do infinitivo) pelo sufixo “-ndo”. O comportamento apresentado, maior discrepância entre média de *heads* na *head* 5 é compartilhado entre diferentes verbos com a mesma flexão verbal.

A Figura 5.12 mostra dois exemplos diferentes. No primeiro, Figura 5.12a, observa-se mais uma vez que a substituição do sufixo de infinitivo (“-r”) pelo sufixo de gerúndio (“-ndo”) resulta em organizações similares das *heads* de atenção. O segundo exemplo, mostrado na Figura 5.12b, onde comparam-se verbos diferentes (“jogando” e “andando”), porém com o mesmo sufixo, as diferenças médias entre *heads* correspondentes são próximas de zero.

Esse padrão, onde as diferenças dos pesos de *heads* correspondentes entre verbos com mesmo sufixo são pequenas e a diferença da *head* 5 é alta para o mesmo verbo com sufixos

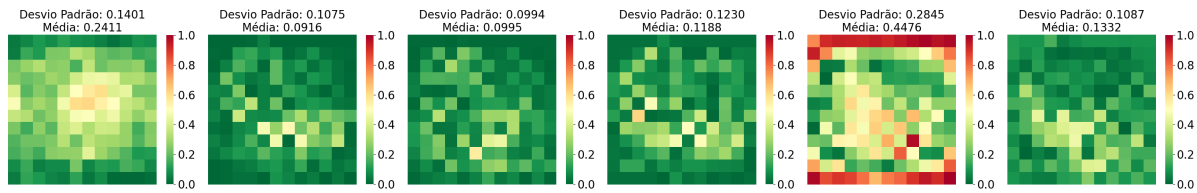
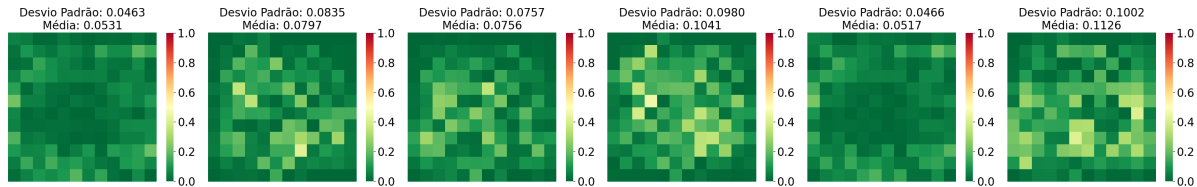
(a) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “jogando” e “jogar”.(b) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “jogando” e “andando”.

Figura 5.12: Comparações entre diferenças de pesos de atenção retornados pelo Modelo-PT. Acima, o mesmo verbo com sufixos diferentes. Embaixo, dois verbos diferentes, porém com o mesmo sufixo.

diferentes, se repete para outros verbos. A Figura 5.13 mostra a repetição desse padrão, dessa vez mostrando verbos diferentes. Primeiro, na Figura 5.13a, a mesma repetição de padrão na *head* 5 quando compara-se a diferença entre “trabalhando” e “trabalhar”, enquanto que a segunda imagem, Figura 5.13b, mostra valores baixos de diferenças, porém, desta vez, com dois verbos compartilhando o mesmo sufixo de infinitivo (“-r”): “correr” e “olhar”.

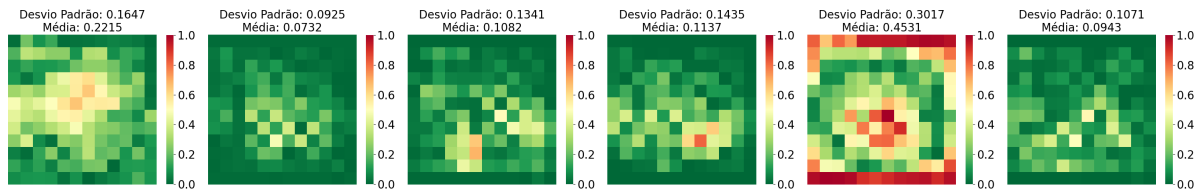
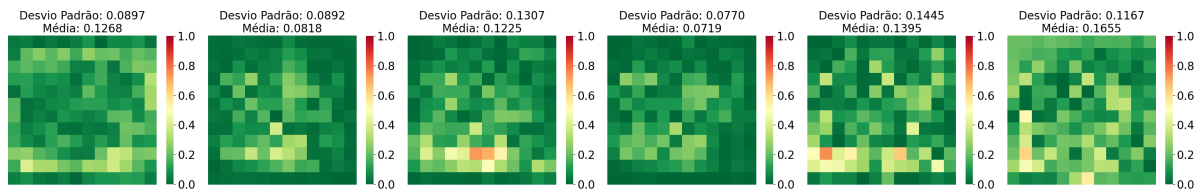
(a) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “trabalhando” e “trabalhar”.(b) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “correr” e “olhar”.

Figura 5.13: Comparações entre diferenças de pesos de atenção retornados pelo Modelo-PT. Acima, o mesmo verbo com sufixos diferentes. Embaixo, dois verbos diferentes, porém com o mesmo sufixo.

Os resultados apresentados até então corroboram com o que os autores em (CLARK et

al., 2019) argumentam em sua conclusão: há uma quantidade substancial de informação linguística contida nos pesos retornados por mecanismos de atenção. Os primeiros casos vistos apontam para uma organização interna entre as *heads* do *Transformer* indicando flexão verbal por meio de sufixos. Para avaliar a influência de outros morfemas gramaticais na organização de mecanismos de atenção, avaliou-se as diferenças entre adjetivos de cores quando ocorro flexão de gênero por meio da mudança dos sufixos “-a” e “-o”.

A escolha por adjetivos de cor se dá pelos resultados apresentados em (PLUMMER et al., 2015), onde os autores inspecionaram as sentenças contidas no conjunto de dados, apontando que adjetivos relativos a cores são os mais frequentes. Após o adjetivo mais utilizado (“*young*”, tradução: jovem), os mais comuns são, em ordem, “*white*”, “*black*”, “*blue*” e “*red*”, que podem ser respectivamente traduzidos para: “branca(o)”, “preta(o)”, “azul” e “vermelha(o)”.

As Figuras 5.14 e 5.15 mostram a influência da flexão de gênero com a diferença das médias dos pesos retornados durante a predição das palavras “branca” e “branco” e “vermelha” e “vermelho” respectivamente. Na primeira comparação, percebe-se uma discrepância maior nas *heads* 4 e 6, enquanto que na segunda imagem, apenas o comportamento da *head* 6 é o mesmo.

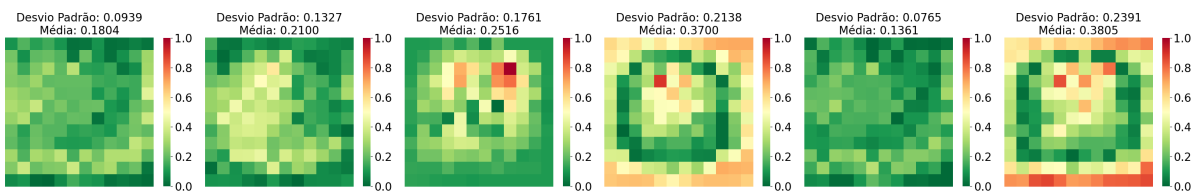


Figura 5.14: Diferenças entre as *heads* do Modelo-PT ao predizer as palavras “branca” e “branco”.

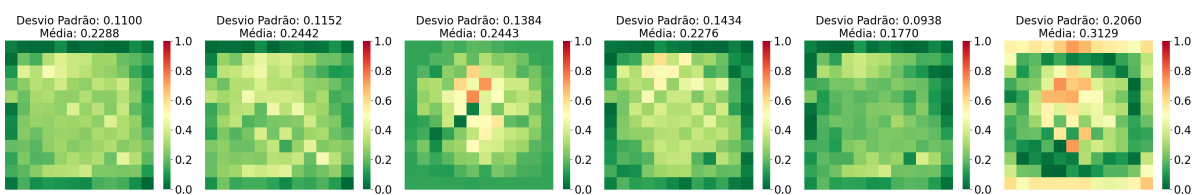


Figura 5.15: Diferenças entre as *heads* do Modelo-PT ao predizer as palavras “vermelha” e “vermelho”.

Apesar dessa disparidade, o funcionamento da *head* 6 se repete para as diferenças entre as palavras “preto” e “preta” e, até mesmo, “branca” e “vermelho”, como é mostrado na Figura 5.16.

Passando para a análise da variação de gênero em substantivos, além da flexão, a heteronímia é outra maneira em que a diferença de gênero pode se manifestar entre seres animados, onde os sexos são apontados com o uso de palavras diferentes como recurso

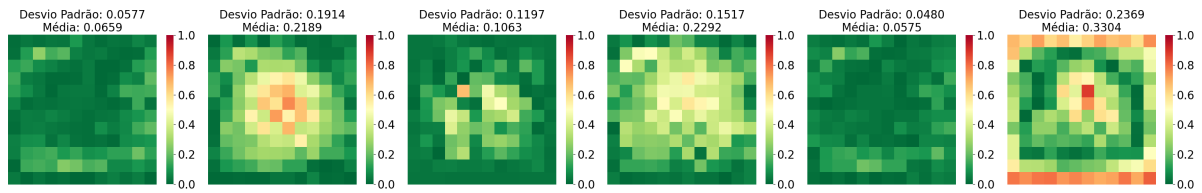
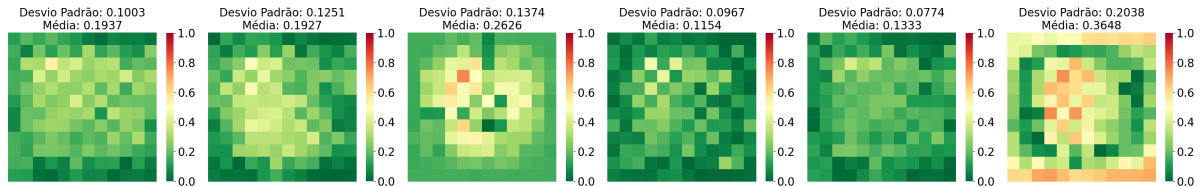
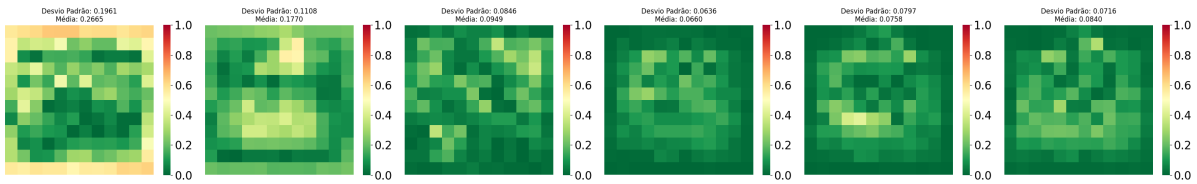
(a) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “preto” e “preta”.(b) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “branca” e “vermelho”.

Figura 5.16: Comparações entre diferenças de pesos de atenção retornados pelo Modelo-PT. Acima, o mesmo adjetivo com sufixos diferentes. Embaixo, dois adjetivos diferentes também com sufixos diferentes.

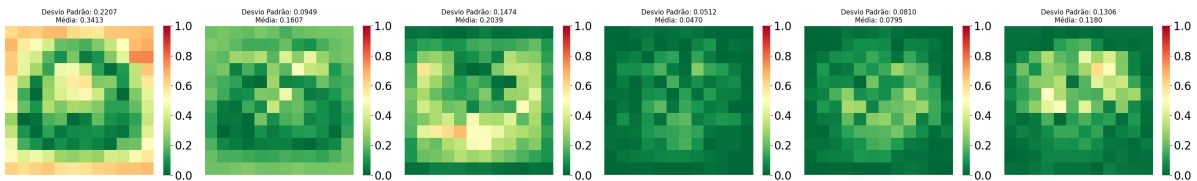
(BECHARA, 2009). Um exemplo de heteronímia presente no conjunto de dados são as palavras “homem” e “mulher” e a diferença entre os pesos de atenção para essas duas palavras é mostrada na Figura 5.17.

Figura 5.17: Diferenças entre as *heads* do Modelo-PT ao prever as palavras “homem” e “mulher”.

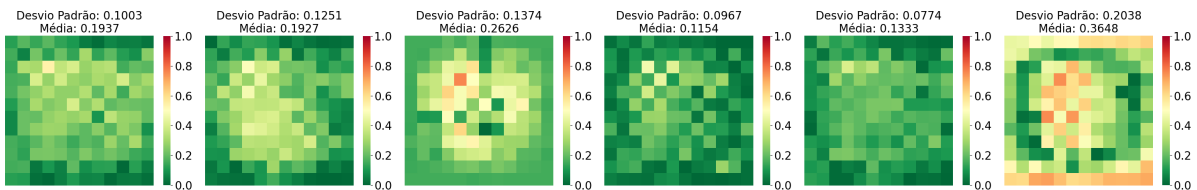
Ainda que de maneira distinta do que foi visto anteriormente, os valores maiores de diferença média na *head* 1 se repetem para outros casos de substantivos representando seres de sexos diferentes, um comportamento que pode revelar um outro padrão sendo seguido para este caso de diferença de gênero das palavras previstas. A Figura 5.18 mostra as diferenças de “foco” das *heads* para substantivos de gêneros opostos e similares.



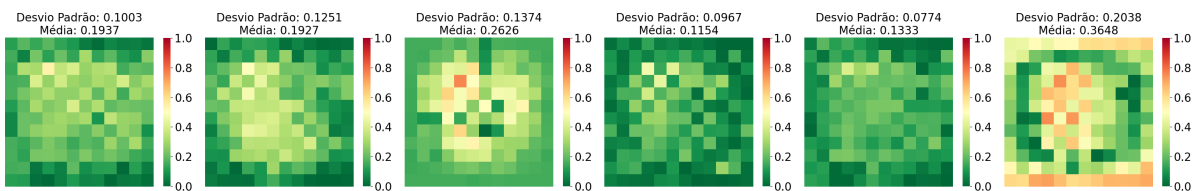
(a) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “homem” e “criança”.



(b) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “mulher” e “cão”.



(c) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “mulher” e “criança”.



(d) Diferenças entre as *heads* do Modelo-PT ao prever as palavras “homem” e “cão”.

Figura 5.18: Comparações entre diferenças de pesos de atenção retornados pelo Modelo-PT.

5.3.7 Experimento 4

O modelo apresentado na Seção 5.2 foi treinado de maneira similar ao que foi feito no Experimento 3, com a adição do *Decoder POS* com os mesmos hiper-parâmetros do *Decoder*, porém com 4 *heads* uma vez que a tarefa de prever POS é considerada mais fácil por haver uma variedade menor de possíveis saídas, uma vez que o “vocabulário” de anotações gramaticais calculado após a extração feita na seção 5.1.2 possui apenas 16 classes diferentes. A escolha de repetir os hiper-parâmetros usados no Experimento 3 para o *Encoder* e o *Decoder* deste experimento se deve para que seja possível uma comparação das sentenças geradas por modelos que fossem em parte similares, com a adição de uma cada de informações a mais (*Decoder POS*).

5.3.8 Resultados 4

Após o treinamento, as mesmas 1000 imagens usadas nos resultados do Experimento 3, na seção 5.3.5, tiveram suas sentenças geradas. A Tabela 5.6 mostra os resultados obtidos com as métricas automáticas (BLEU 1-4 e METEOR) para os modelos Modelo-LIC-EN (treinado com Flickr30k-EN) e Modelo-LIC-PT (treinado com Flickr30k-PT).

Model	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR
Modelo-LIC-EN	60.50	42.63	29.26	20.05	18.98
Modelo-LIC-PT	53.38	34.62	22.43	14.48	17.74

Tabela 5.6: Métricas obtidas para cada modelo treinado e avaliado com seus respectivos conjuntos de dados utilizando o modelo proposto (Modelo-LIC).

Os valores obtidos para as métricas foram menores do que os do *Transformer* treinado para o Experimento 3. Apesar disso, as sentenças preditas pelo Modelo-LIC-PT têm um tamanho médio maior, $12,11 \pm 4,11$, quando comparadas com as geradas pelo Modelo-PT, $10,60 \pm 2,83$, o que pode ser um indicativo de sentenças mais informativas. As Figuras 5.19, 5.20, 5.21, 5.22 e 5.23 mostram exemplos, escolhidos de forma aleatória, de sentenças geradas com os modelos Modelo-PT e Modelo-LIC-PT. Nas figuras, pode-se observar que o Modelo-LIC-PT produz sentenças maiores e mais informativas, mesmo que as vezes errando ao adicionar informações repetidas (Figura 5.19) ou errando cores de objetos (Figura 5.20).



Figura 5.19:

- **Modelo-PT:** um rapaz com uma camisa vermelha está a saltar
- **Modelo-LIC-PT:** um rapaz com uma camisa vermelha e um homem com uma camisa vermelha



Figura 5.20:

- **Modelo-PT:** uma mulher está a conduzir um autocarro vermelho
- **Modelo-LIC-PT:** uma mulher em uma camisa rosa está em frente a um ônibus vermelho com uma camisa vermelha



Figura 5.21:

- **Modelo-PT:** um grupo de pessoas está posando para uma imagem
- **Modelo-LIC-PT:** um grupo de pessoas está de pé em uma bicicleta



Figura 5.22:

- **Modelo-PT:** um homem e uma mulher estão a falar com uma mulher
- **Modelo-LIC-PT:** duas mulheres e um homem estão sentados em uma mesa com um laptop



Figura 5.23:

- **Modelo-PT:** um grupo de pessoas está a falar com uma mulher
- **Modelo-LIC-PT:** um grupo de pessoas está a olhar para algo em uma mesa

Assim, observa-se que o Modelo-LIC-PT produz descrições menos genéricas quando comparadas às mesmas sentenças geradas pelo Modelo-PT, dando indícios de que o uso de classes gramaticais durante o treinamento contribui para a geração de palavras e seu posicionamento na sentença ao adicionar conhecimento linguístico presente nos dados de treino.

5.4 RESUMO DO CAPÍTULO

Este capítulo abordou os experimentos realizados para endereçar as perguntas criadas na seção 2.4. Inicialmente abordou-se a metodologia empregada, em que os conjuntos utilizados foram apresentados, bem como a configuração das classes gramaticais após serem extraídas com a biblioteca *spaCy* e como foram realizadas as avaliações quantitativa e qualitativa. Em seguida a arquitetura do modelo treinado, utilizando classes gramaticais, foi apresentada junto com a justificativa do uso de informações sintáticas e como estas foram empregadas durante o treinamento. Finalizando com os experimentos feitos ao longo do trabalho, validando o uso de conjuntos de dados traduzidos e explorando erros encontrados, para, em seguida, explorar os efeitos de mudanças morfológicas em mecanismos de atenção para a tarefa de *Image Captioning*, culminando com o uso de classes gramaticais no processo de treinamento da arquitetura de rede neural proposta.

CONCLUSÃO E TRABALHOS FUTUROS

Este capítulo apresenta a conclusão do trabalho feito. Primeiramente é apresentada uma visão geral da temática de *Image Captioning*, juntamente a motivação para investigação e os problemas apontados que levaram ao objetivo principal trabalho. Em seguida são listadas as contribuições dos experimentos feitos e trabalhos publicados. Por fim são apresentadas as limitações dos experimentos feitos e direcionamentos para trabalhos futuros.

A tarefa de *Image Captioning* busca automatizar uma tarefa que é natural ao ser humano: visualizar e compreender o conteúdo de uma imagem para, em seguida, descrevê-la com uma sentença em linguagem natural e semanticamente correta. Pela descrição da tarefa percebe-se que os sistemas propostos permeiam áreas distintas da Inteligência Artificial: Visão Computacional e Processamento de Linguagem Natural. Sistemas de descrição de imagens de forma automatizada têm aplicações que vão desde melhorias em acessibilidade web até o auxílio de sistemas de interação humano computador, sendo uma área de pesquisa relativamente recente, com primeiros estudos reportados há menos de 20 anos atrás.

Mais da metade dos artigos endereçados para a tarefa de *Image Captioning* são idealizados para a língua inglesa, que conta com pouca variação de palavras e, geralmente, com uma ordem engessada de palavras, portanto não sendo representativa o suficiente para estudos na área com diferentes linguagens. Investigar as predições de um modelo gerado para o inglês e categorizar os erros existentes pode não ser o suficiente para alcançar melhorias em outros idiomas, o que leva a primeira motivação para este trabalho: a falta de representatividade da língua portuguesa em estudos na tarefa de *Image Captioning*. Ademais, estudos mostram que o uso de conhecimento linguístico durante a modelagem de sistemas de PLN aumenta o poder preditivo desses sistemas, tornando-os mais robustos para línguas que podem variar bastante em suas propriedades.

Este trabalho concentra-se na área de PLN para analisar o uso de informações provenientes da linguagem humana e que estão presentes nos dados de treinamento (e nas sentenças geradas) usados para aprendizados de sistemas de *Image Captioning*. Os experimentos foram realizados com conjuntos de dados traduzidos para o Português, portanto,

para avaliar as possibilidades do uso de dados traduzidos automaticamente, no primeiro experimento comparou-se as métricas obtidas ao treinar o mesmo modelo de rede neural com dois conjuntos de dados diferentes, um com sentenças originais em inglês e outro traduzido de maneira automática. Os resultados mostraram pouca diferença em relação às métricas calculadas, sendo feita em seguida uma avaliação qualitativa utilizando anotadores externos para marcar a presença de erros categorizados pelos autores.

O experimento 3 teve como objetivo avaliar o impacto de flexões nominais e verbais, bem como heteronímias, nos pesos de atenção retornados por mecanismos de *Multi-Head Attention* de um *Transformer* treinado com o conjunto de dados *Flickr30k*. Ao inspecionar as diferenças das médias de atenção dada para palavras com morfemas gramaticais distintos e similares, notou-se um padrão de comportamento entre as *heads*, indicando que informações morfológicas são aprendidas durante o treinamento dos modelos. Por último, avaliou-se o uso de classes gramaticais no processo de treinamento da arquitetura proposta na seção 5.2.

6.1 PUBLICAÇÕES

Os resultados 1 (5.3.2) e 2 (5.3.4) deste trabalho foram publicados no 24th *International Conference on Enterprise Information Systems (ICEIS 2022)* (GONDIM.; CLARO.; SOUZA., 2022). Para atingir o Objetivo Específico 1, foram apresentados os Resultados 1 e 2 com as considerações a respeito da comparação de métricas entre modelos similares treinados com conjuntos de dados de diferentes linguagens: *Flickr8k* original e traduzido para o português. Com o auxílio de anotadores humanos, a presença de erros comuns em um modelo de descrição de imagens foi contabilizada para 100 imagens. As questões levantadas após analisar os erros buscam explorar as falhas existentes com o uso de classes gramaticais para guiar o treinamento de um modelo com arquitetura inspirada em um *Transformer*.

No experimento 3, atrelado ao Objetivo Específico 2, foi feita uma análise do mecanismo de *Multi-Head Attention* de uma rede *Transformer* para a tarefa de *Image Captioning*, onde comparou-se o comportamento apresentado por cada *head* com pares de palavras com mudanças morfológicas de gênero e modo verbal. O estudo segue a análise de padrões de comportamento de *heads* de atenção, abordada em outros trabalhos, mas com o diferencial de ser feito com entradas e saídas multimodais (imagem e texto respectivamente) e pretende-se realizar a submissão para o periódico *Natural Language Engineering*.

Já no experimento 4, também ligado ao Objetivo Específico 2, foi proposta uma arquitetura de rede neural para a tarefa de descrição automatizada de imagens. Durante o treinamento dessa rede, as classes gramaticais de cada palavra das sentenças de treino são inseridas juntamente com as sentenças e as imagens. Ainda que não tenha se observado melhoras nas métricas automatizadas, as sentenças geradas pelo modelo proposto têm tamanho médio maior, indicando inferências mais informativas que foram observadas quando comparadas com as sentenças geradas por um modelo similar sem o uso de classes gramaticais. Pretende-se publicar essas descobertas também no periódico *Natural Language Engineering*.

6.2 LIMITAÇÕES

1. **Erros de tradução:** a primeira limitação constatada durante a sequencia do trabalho foram os erros de tradução automática cometidos pela ferramenta utilizada (LibreTranslate¹). Como esses erros são inseridos no processo de treinamento, a performance de modelos iguais treinados com conjuntos de dados diferentes tende a ser melhor para as sentenças originais em inglês.
2. **Baixa variabilidade das sentenças de teste:** o conjunto de dados *Flickr30k* possui 31014 imagens disponibilizadas para treinamento, validação e teste, cada uma com 5 sentenças, totalizando 155070 sentenças. O número é alto, mas as imagens possuem muitas características em comum entre si, dificultando a análise de outros possíveis casos. No trabalho (PLUMMER et al., 2015), os autores apontam que a palavra *man* é a que aparece mais vezes entre substantivos, sendo seguida por *woman*, mas com metade da quantidade de aparições, enquanto que o terceiro substantivo que mais aparece é *people*, mais uma vez caindo pela metade em aparições quando comparado com *woman*.
3. **Falhas na anotação do conjunto de dados:** apesar de ser um conjunto de dados utilizado na literatura, o *Flickr30k* possui sentenças que não fazem sentido ou que foram anotadas com opiniões pessoais pelos anotadores. Algumas das frases encontradas são: “I don’t see a picture i don’t see a picture i don’t see a picture i don’t see a picture”, “You know i am looking like Justin Bieber” e “What is this I don’t even”. Além de representarem uma sentença a menos no treinamento, esses erros atrapalham o aprendizado do modelo.
4. **Viéses na coleta de imagens e dos anotadores:** Como visto na segunda limitação listada, a palavra “woman” aparece metade das vezes quando comparada com “man”. Além de diminuir a representatividade, podendo influenciar em predições piores quando há mulheres na imagem, uma das sentenças é misógina: “Why aren’t these girls in the kitchen instead of playing with toys?”.

6.3 TRABALHOS FUTUROS

Ainda que o uso de classes gramaticais não tenha resultado em aumento das métricas utilizadas, o uso de POS em outros trabalhos, principalmente em (WANG et al., 2022), além de trabalhos advogando pelo uso de conhecimento linguístico em modelos de PLN, como (BENDER, 2009, 2011, 2013), sustentam abordagens que empreguem conhecimento linguístico das sentenças presentes nos conjuntos de dados. O experimento 3 mostrou que há um aprendizado de flexões e heteronímias durante o treinamento e os autores em (CLARK et al., 2019) mostram que as *heads* de atenção se organizam para indicar classes gramaticais, portanto, uma vez que essas informações são aprendidas, bem como a associação de sentenças dada uma imagem, aproveitar que elas podem ser extraídas

¹<https://github.com/LibreTranslate/LibreTranslate>

e utilizadas mostra-se como um caminho a ser seguido em busca de sistemas de *Image Captioning* para o Português.

Além de classes gramaticais, facilmente extraídas com o pacote *spaCy*, pretende-se explorar também relações retornadas por analisadores de dependência, que podem ser utilizadas de forma análoga ao *Positional Encoding* de um *Transformer*, porém estabelecendo uma ordem de relação entre cada uma das palavras da sentença. Outro caminho a ser explorado é o uso de lematização, processo em que uma palavra é transformada em seu radical, de forma similar ao *Decoder POS*, adicionando mais uma camada de sumarização da imagem com uma sentença retornada mais próxima da sentença alvo.

REFERÊNCIAS BIBLIOGRÁFICAS

AGHDAM, H. H.; HERAVI, E. J. *Guide to Convolutional Neural Networks: A Practical Application to Traffic-Sign Detection and Classification*. 1st. ed. [S.l.]: Springer Publishing Company, Incorporated, 2017. ISBN 331957549X.

AKER, A.; GAIZAUSKAS, R. Generating image descriptions using dependency relational patterns. In: *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*. USA: Association for Computational Linguistics, 2010. (ACL '10), p. 1250–1258.

ANDERSON, P. et al. *SPICE: Semantic Propositional Image Caption Evaluation*. arXiv, 2016. Disponível em: <<https://arxiv.org/abs/1607.08822>>.

BAHDANAU, D.; CHO, K.; BENGIO, Y. *Neural Machine Translation by Jointly Learning to Align and Translate*. arXiv, 2014. Disponível em: <<https://arxiv.org/abs/1409.0473>>.

BAI, S.; AN, S. A survey on automatic image caption generation. *Neurocomputing*, v. 311, p. 291–304, 2018. ISSN 0925-2312. Disponível em: <<https://www.sciencedirect.com/science/article/pii/S0925231218306659>>.

BANERJEE, S.; LAVIE, A. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In: *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*. Ann Arbor, Michigan: Association for Computational Linguistics, 2005. p. 65–72. Disponível em: <<https://aclanthology.org/W05-0909>>.

BECHARA, E. *Moderna gramática portuguesa*. [S.l.]: Editora Lucerna, 2009. ISBN 978-85-209-3049-6.

BENDER, E. English isn't generic for language, despite what nlp papers might lead you to believe. In: *Symposium and Data Science and Statistics*. [s.n.], 2019. [Online; accessed 15-may-2020]. Disponível em: <<http://faculty.washington.edu/ebender/papers/Bender-SDSS-2019.pdf>>.

BENDER, E. M. Linguistically naïve != language independent: Why NLP needs linguistic typology. In: *Proceedings of the EACL 2009 Workshop on the Interaction between Linguistics and Computational Linguistics: Virtuous, Vicious or Vacuous?* Athens, Greece: Association for Computational Linguistics, 2009. p. 26–32. Disponível em: <<https://www.aclweb.org/anthology/W09-0106>>.

BENDER, E. M. On achieving and evaluating language-independence in nlp. *Linguistic Issues in Language Technology*, v. 6, Oct. 2011. Disponível em: <<https://journals.colorado.edu/index.php/lilt/article/view/1239>>.

BENDER, E. M. Morphology: Introduction. In: _____. *Linguistic Fundamentals for Natural Language Processing: 100 Essentials from Morphology and Syntax*. Cham: Springer International Publishing, 2013. ISBN 978-3-031-02150-3. Disponível em: <https://doi.org/10.1007/978-3-031-02150-3_2>.

BERNARDI, R. et al. *Automatic Description Generation from Images: A Survey of Models, Datasets, and Evaluation Measures*. 2017.

BIRD, S.; KLEIN, E.; LOPER, E. *Natural Language Processing with Python*. [S.l.: s.n.], 2009. ISBN 978-0-596-51649-9.

BROWNLEE, J. *Deep Learning for Natural Language Processing: Develop Deep Learning Models for your Natural Language Problems*. Machine Learning Mastery, 2017. Disponível em: <https://books.google.com.br/books?id=_pmoDwAAQBAJ>.

CARREIRA, J. et al. Semantic segmentation with second-order pooling. In: . [S.l.: s.n.], 2012. ISBN 978-3-642-33785-7.

CHANG, X. et al. Scene graphs: A survey of generations and applications. *CoRR*, abs/2104.01111, 2021. Disponível em: <<https://arxiv.org/abs/2104.01111>>.

CHO, K. et al. *On the Properties of Neural Machine Translation: Encoder-Decoder Approaches*. arXiv, 2014. Disponível em: <<https://arxiv.org/abs/1409.1259>>.

CHO, K. et al. *Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation*. arXiv, 2014. Disponível em: <<https://arxiv.org/abs/1406.1078>>.

CHOLLET, F. *Deep Learning with Python*. 1st. ed. USA: Manning Publications Co., 2017. ISBN 1617294438.

CLARK, K. et al. What does BERT look at? an analysis of BERT's attention. In: *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Florence, Italy: Association for Computational Linguistics, 2019. p. 276–286. Disponível em: <<https://aclanthology.org/W19-4828>>.

CORMEN, T. H. et al. *Introduction to Algorithms, Third Edition*. 3rd. ed. [S.l.]: The MIT Press, 2009. ISBN 0262033844.

CORNIA, M. et al. Meshed-Memory Transformer for Image Captioning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. [S.l.: s.n.], 2020.

CUNHA, C.; CINTRA, L. *Nova gramática do português contemporâneo*. 7 ed. ed. Rio de Janeiro :: Lexicon, 2016.

- DOGNIN, P. L. et al. *Adversarial Semantic Alignment for Improved Image Captions*. arXiv, 2018. Disponível em: <<https://arxiv.org/abs/1805.00063>>.
- FARHADI, A. et al. Every picture tells a story: Generating sentences from images. In: *Proceedings of the European Conference on Computer Vision (ECCV'10)*. [S.l.: s.n.], 2010. v. 6314, p. 15–29. ISBN 978-3-642-15560-4.
- GATT, A.; KRAHMER, E. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *J. Artif. Int. Res.*, AI Access Foundation, El Segundo, CA, USA, v. 61, n. 1, p. 65–170, jan 2018. ISSN 1076-9757.
- GONDIM., J.; CLARO., D.; SOUZA., M. Towards image captioning for the portuguese language: Evaluation on a translated dataset. In: INSTICC. *Proceedings of the 24th International Conference on Enterprise Information Systems - Volume 1: ICEIS*. SciTePress, 2022. p. 384–393. ISBN 978-989-758-569-2. ISSN 2184-4992. Disponível em: <<https://doi.org/10.5220/0011080000003179>>.
- GOODFELLOW, I.; BENGIO, Y.; COURVILLE, A. *Deep Learning*. [S.l.]: MIT Press, 2016. <<http://www.deeplearningbook.org>>.
- GRAVES, A. *Supervised Sequence Labelling with Recurrent Neural Networks*. [S.l.: s.n.], 2012. ISBN 978-3-642-24796-5.
- HE, K. et al. *Deep Residual Learning for Image Recognition*. arXiv, 2015. Disponível em: <<https://arxiv.org/abs/1512.03385>>.
- HE, X. et al. Image caption generation with part of speech guidance. *Pattern Recognition Letters*, v. 119, 10 2017.
- HOCHREITER, S.; SCHMIDHUBER, J. Long short-term memory. *Neural computation*, v. 9, p. 1735–80, 12 1997.
- HODOSH, M.; YOUNG, P.; HOCKENMAIER, J. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, v. 47, p. 853–899, 05 2013.
- HUANG, L. et al. *Attention on Attention for Image Captioning*. arXiv, 2019. Disponível em: <<https://arxiv.org/abs/1908.06954>>.
- HUGHES, J. F. et al. *Computer Graphics: Principles and Practice*. 3. ed. Upper Saddle River, NJ: Addison-Wesley, 2013. ISBN 978-0-321-39952-6.
- JURAFSKY, D.; MARTIN, J. H. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. [S.l.: s.n.], 2022. <<https://web.stanford.edu/~jurafsky/slp3/>>. Acessado em: 20-02-2022.
- KARPATHY, A.; FEI-FEI, L. *Deep Visual-Semantic Alignments for Generating Image Descriptions*. arXiv, 2014. Disponível em: <<https://arxiv.org/abs/1412.2306>>.

KARPATHY, A.; JOULIN, A.; FEI-FEI, L. *Deep Fragment Embeddings for Bidirectional Image Sentence Mapping*. arXiv, 2014. Disponível em: <<https://arxiv.org/abs/1406.5679>>.

KLEIN, G. et al. OpenNMT: Open-source toolkit for neural machine translation. In: *Proceedings of ACL 2017, System Demonstrations*. Vancouver, Canada: Association for Computational Linguistics, 2017. p. 67–72. Disponível em: <<https://www.aclweb.org/anthology/P17-4012>>.

LIMA, E. L. *Álgebra Linear*. 1. ed. [S.l.]: IMPA, 2014. ISBN 978-85-244-0390-3.

LIN, T.-Y. et al. *Microsoft COCO: Common Objects in Context*. 2015.

LU, J. et al. *Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning*. 2017.

MANNING, C. D.; SCHÜTZE, H. *Foundations of Statistical Natural Language Processing*. Cambridge, Massachusetts: The MIT Press, 1999. Disponível em: <<http://nlp.stanford.edu/fsnlp/>>.

MIELKE, S. J. *Language diversity in ACL 2004 - 2016*. 2016. Disponível em: <<https://sjmielke.com/acl-language-diversity.htm>>.

MIKOLOV, T. et al. Recurrent neural network based language model. In: . [S.l.: s.n.], 2010. v. 2, p. 1045–1048.

MILTENBURG, E. van; ELLIOTT, D. *Room for improvement in automatic image description: an error analysis*. arXiv, 2017. Disponível em: <<https://arxiv.org/abs/1704.04198>>.

MINAEE, S. et al. Image segmentation using deep learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, p. 1–1, 2021.

MONTANI, I. et al. *explosion/spaCy: v3.3.0: Improved speed, new trainable lemmatizer, and pipelines for Finnish, Korean and Swedish*. Zenodo, 2022. Disponível em: <<https://doi.org/10.5281/zenodo.6504092>>.

ORDONEZ, V.; KULKARNI, G.; BERG, T. L. Im2text: Describing images using 1 million captioned photographs. In: *Neural Information Processing Systems (NIPS)*. [S.l.: s.n.], 2011.

PAN, J.-Y. et al. Automatic image captioning. In: . [S.l.: s.n.], 2004. v. 3, p. 1987 – 1990 Vol.3. ISBN 0-7803-8603-5.

PAN, Y. et al. *X-Linear Attention Networks for Image Captioning*. arXiv, 2020. Disponível em: <<https://arxiv.org/abs/2003.14080>>.

PAPINENI, K. et al. Bleu: A method for automatic evaluation of machine translation. In: *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*. USA: Association for Computational Linguistics, 2002. (ACL '02), p. 311–318. Disponível em: <<https://doi.org/10.3115/1073083.1073135>>.

PENNINGTON, J.; SOCHER, R.; MANNING, C. D. Glove: Global vectors for word representation. In: *Empirical Methods in Natural Language Processing (EMNLP)*. [s.n.], 2014. p. 1532–1543. Disponível em: <<http://www.aclweb.org/anthology/D14-1162>>.

PLUMMER, B. A. et al. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In: *2015 IEEE International Conference on Computer Vision (ICCV)*. [S.l.: s.n.], 2015. p. 2641–2649.

RADEMAKER, A. et al. Universal dependencies for portuguese. In: *Proceedings of the Fourth International Conference on Dependency Linguistics (Depling)*. Pisa, Italy: [s.n.], 2017. p. 197–206. Disponível em: <<http://aclweb.org/anthology/W17-6523>>.

ROBERTSON, S. Understanding inverse document frequency: On theoretical arguments for idf. *Journal of Documentation - J DOC*, v. 60, p. 503–520, 10 2004.

ROSA, G. M. et al. *A cost-benefit analysis of cross-lingual transfer methods*. arXiv, 2021. Disponível em: <<https://arxiv.org/abs/2105.06813>>.

RUSSELL, S. J.; NORVIG, P. *Artificial intelligence : a modern approach*. 4th edition. ed. Boston: Pearson, 2020. ISBN 0134610997; 9780134610993.

SANTOS, G. O. dos; COLOMBINI, E. L.; AVILA, S. #pracegover: A large dataset for image captioning in portuguese. *Data*, v. 7, n. 2, 2022. ISSN 2306-5729. Disponível em: <<https://www.mdpi.com/2306-5729/7/2/13>>.

SHARIF, N. et al. Vision to language: Methods, metrics and datasets. In: _____. [S.l.: s.n.], 2020. p. 9–62. ISBN 978-3-030-49723-1.

SILVA, P. N. d. *Manual de introdução aos estudos linguísticos*. [s.n.], 2010. Disponível em: <<http://hdl.handle.net/10400.2/11900>>.

STEFANINI, M. et al. *From Show to Tell: A Survey on Deep Learning-based Image Captioning*. arXiv, 2021. Disponível em: <<https://arxiv.org/abs/2107.06912>>.

SULTANA, F.; SUFIAN, A.; DUTTA, P. Advancements in image classification using convolutional neural network. *CoRR*, abs/1905.03288, 2019. Disponível em: <<http://arxiv.org/abs/1905.03288>>.

TAN, M.; LE, Q. V. *EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks*. 2020.

TSENG, B.-H. et al. *A Generative Model for Joint Natural Language Understanding and Generation*. 2020.

- VASWANI, A. et al. Attention is all you need. In: GUYON, I. et al. (Ed.). *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2017. v. 30. Disponível em: <<https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>>.
- VEDANTAM, R.; ZITNICK, C. L.; PARIKH, D. *CIDEr: Consensus-based Image Description Evaluation*. arXiv, 2014. Disponível em: <<https://arxiv.org/abs/1411.5726>>.
- VIG, J. *A Multiscale Visualization of Attention in the Transformer Model*. arXiv, 2019. Disponível em: <<https://arxiv.org/abs/1906.05714>>.
- VINYALS, O. et al. *Show and Tell: A Neural Image Caption Generator*. 2015.
- WANG, A. et al. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In: *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Brussels, Belgium: Association for Computational Linguistics, 2018. p. 353–355. Disponível em: <<https://aclanthology.org/W18-5446>>.
- WANG, D. et al. Separate syntax and semantics: Part-of-speech-guided transformer for image captioning. *Applied Sciences*, v. 12, n. 23, 2022. ISSN 2076-3417. Disponível em: <<https://www.mdpi.com/2076-3417/12/23/11875>>.
- Web Accessibility Initiative. *Introduction to Web Accessibility*. 2022. <<https://www.w3.org/WAI/fundamentals/accessibility-intro/#context>>. [Online; acessado em 3-Abril-2022].
- XU, K. et al. *Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*. 2016.
- YANG, Y. et al. Corpus-guided sentence generation of natural images. In: *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*. Edinburgh, Scotland, UK.: Association for Computational Linguistics, 2011. p. 444–454. Disponível em: <<https://aclanthology.org/D11-1041>>.
- YAO, B. Z. et al. I2t: Image parsing to text description. *Proceedings of the IEEE*, v. 98, n. 8, p. 1485–1508, 2010.
- YOUNG, P. et al. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *TACL*, v. 2, p. 67–78, 12 2014.
- ZHANG, A. et al. *Dive into Deep Learning*. arXiv, 2021. Disponível em: <<https://arxiv.org/abs/2106.11342>>.
- ZHANG, J. et al. Integrating part of speech guidance for image captioning. *IEEE Transactions on Multimedia*, v. 23, p. 92–104, 2021.