

CONSTRUÇÃO DE UM CORPUS PARALELO ALINHADO PARA A TRADUÇÃO AUTOMÁTICA ESTATÍSTICA

Wilker Ferreira Aziz (ICMC-USP)

Thiago Alexandre Salgueiro Pardo (ICMC-USP)

Ivandr  Paraboni (EACH-USP)

A necessidade da tradu  o de grandes quantidades de textos em pouco tempo estimulou, no passado, a pesquisa de m todos de Tradu  o Autom tica (TA). Atualmente, a Internet aprimorou o desafio devido a enorme quantidade de informa  o multil ng e. A tradu  o entre pares de l nguas distantes incentivou a busca por modelos estat sticos de TA, pois pouco conhecimento ling  stico   necess rio nessa abordagem. C rpus paralelos s o a base para o desenvolvimento desses modelos, pois, a partir deles   extra do automaticamente todo o conhecimento necess rio para a tradu  o. Por c rpus paralelo, entende-se um conjunto de textos em uma l ngua-fonte acompanhados das tradu  es na l ngua-alvo. Esses c rpus podem ser alinhados lexicalmente e/ou sentencialmente, ou seja, s o indicadas, em cada texto, as correspond ncias entre suas palavras e/ou senten as com o texto na outra l ngua. Existem v rios m todos para se produzir automaticamente o alinhamento entre textos paralelos, classificados como emp ricos, ling  sticos ou h bridos, de acordo com o n vel de conhecimento ling  stico que utilizam. Os m todos emp ricos baseiam-se em medidas do texto, tais como n mero de palavras, quantidade de caracteres, ocorr ncia de caracteres especiais (letras mai sculas e pontua  o), similaridades e padr es, dispensando, assim, qualquer conhecimento ling  stico; os m todos ling  sticos, por sua vez, utilizam conhecimento sobre as l nguas envolvidas, tais como gloss rios, l xicos, listas de elementos que freq entemente s o correspondentes entre as l nguas e regras gramaticais, entre outros recursos; por fim, os m todos h bridos fazem uso de recursos de ambos os m todos anteriores para realizar o alinhamento. H  diversos c rpus paralelos dispon veis para v rias l nguas, por exemplo, o c rpus composto pelas cartas da ONU (Organiza  o das Na  es Unidas), traduzidas para os mais diversos idiomas. Da mesma forma, h  diversas ferramentas de alinhamento textual, como o GIZA++ e o LIHLA, sendo que este  ltimo foi desenvolvido no NILC (N cleo Interinstitucional de Ling  stica Computacional), um dos maiores centros de pesquisa em Processamento de L nguas Naturais no Brasil. Para o desenvolvimento de bons tradutores autom ticos estat sticos, faz-se necess rio um grande c rpus paralelo, em geral, com mais de 200 milh es de palavras, que, idealmente, deve conter alinhamentos corretos. Prop e-se, neste trabalho, a constru  o de um c rpus paralelo alinhado representativo para o par de l nguas portugu s-espanhol, com a finalidade de desenvolver um tradutor autom tico para tais l nguas. Posteriormente, pretende-se estender esse c rpus para a l ngua inglesa. A metodologia para a constru  o desse c rpus  : coleta dos textos e suas tradu  es a partir de fonte confi vel, se poss vel, revisados por humanos; escolha de uma ferramenta de alinhamento sentencial, adequa  o dos textos ao formato de entrada exigido pela ferramenta e gera  o dos alinhamentos sentenciais; produzido o alinhamento sentencial, faz-se necess ria, para garantia de resultados consistentes, a p s-ed  o/revis o dos alinhamentos (manualmente, em geral); escolha de uma ferramenta de alinhamento lexical, adequa  o dos textos alinhados sentencialmente ao formato de entrada exigido pela ferramenta e gera  o dos alinhamentos lexicais; p s-ed  o/revis o dos alinhamentos lexicais. Inicialmente, pretende-se construir o c rpus com artigos da Revista FAPESP (publicada em portugu s, espanhol e ingl s). As ferramentas de alinhamento TCAalign e LIHLA dever o ser utilizadas neste trabalho. Ferramentas visuais podem ser usadas para facilitar o processo de revis o, tal como o VisualTCA, tamb m desenvolvido no NILC.   importante dizer que, em princ pio, o alinhamento lexical n o se faz necess rio para o desenvolvimento de tradutores estat sticos, mas   sabido que este recurso pode aprimorar os resultados obtidos. Com a disponibiliza  o do c rpus paralelo alinhado, h  a possibilidade de se desenvolver tradutores estat sticos propriamente ditos, via o uso de toolkits dispon veis para a comunidade de pesquisa, por exemplo, o GIZA++ e o sistema Rewriter. Os c rpus e os tradutores a serem constru dos s o a grande contribui  o deste trabalho.