

WSD on hierarchical-organized semantics

Emanuele Rucci, 2053183

1 Abstract

2 Word Sense Disambiguation (WSD) poses a
3 significant challenge within the field of Natural
4 Language Processing (NLP), requiring the
5 accurate determination of a word's intended
6 meaning in a given context. This document aims
7 to explore diverse strategies documented in the
8 literature for solving this issue. These strategies
9 are then applied to the provided dataset,
10 evaluating their effectiveness.

11 Additionally, the document analyzes the
12 connection between identifying coarse-grained
13 senses and enhancing the precise classification
14 of fine-grained senses, and vice versa. By
15 investigating this relationship, the document
16 aims to uncover the synergistic effects of these
17 seemingly disparate processes.

18 1 Introduction

19 WSD aims to solve the ambiguity of word meaning
20 in context; when a word has multiple meaning
21 (polysemy) is fundamental to disambiguate the
22 given target word in order to fully understand the
23 complete meaning of a sentence. In order to
24 accomplish this task two big group of technique
25 have been applied in literature during time: fully
26 neural approaches ([1][2] and many others) and
27 knowledge-based approaches [3][4]. Here in this
28 document, we will explore neural approaches
29 applied on the given dataset. Moreover, a couple of
30 experiment have been devoted to study how might
31 be possible to combine in a single model the power
32 to extrapolate in hierarchical way deeper semantic
33 in form of word senses.

34 2 Methodology

35 As said, have been studied models in the neural
36 approach family for WSD. All the models
37 generated consists of an encoder and a
38 classification part [2]. The encoder extracts the

39 embeddings of the sentence and the features are
40 used to predict the senses using the classifier.

41 2.1 Encoder

42 BERT model as the encoder [5]: it is responsible of
43 tokenization and generate the contextualized
44 embedding for the sentence (the hidden states for
45 each token with dimension H).

46 For each word to disambiguate the network
47 takes the mean of the last 4 layer of BERT [6].

48 With these settings the task is being translated in
49 a token-to-token tagging task: assign to each word
50 to be disambiguated the sense predicted. However,
51 an extra step is required to avoid the token
52 fragmentation of a word and retrieve the complete
53 embedding information: the BERT tokenizer might
54 split a single word in 2 or more tokens; when this
55 happens, the model identifies the case and then,
56 considering a list of tokens mapped to the same
57 word, their hidden states are averaged:

$$f = \frac{1}{k} \sum_{l=1}^{k-1} h_{j+l}$$

59 In the equation h_i is the hidden state of token in
60 the list of the k tokens mapped to the same word.

61 Image 1, (a) shows a schema about the encoding
62 process.

63 Moreover, have been tested two implementation
64 strategy for the forward of the linear layer:

- 65 1. f contains the embeddings for all the tokens
66 of the sentence;
- 67 2. f contains only the embeddings of the word
68 to be disambiguated: from the f tensor are
69 removed all the h_i correspondent to token i
70 which are not mapped to a word to be
71 disambiguated.

72 2.2 Classifier

73 All the classifiers are fully connected linear-layers
74 and have been used 3 layers in the different settings
75 of the experiments:

$$76 L_1(x) = W_1 x + b_i \text{ and } W_1 \in \mathbb{R}^{|CG| \times H}$$

77 $L_2(x) = W_2 x + b_i$ and $W_2 \in \mathbb{R}^{|FG| \times H}$
78 $L_3(x) = W_3 x + b_i$ and $W_3 \in \mathbb{R}^{|FG| \times |CG|}$
79 CG and FG represent the sets of coarse and fine
80 grained senses in the dataset.

81 Before to apply any of this linear layer, a mask is
82 applied to the logits (the output of the encoder) in
83 order to ignore all the logits value for the tokens
84 mapped to word that doesn't have to be
85 disambiguated; differently from [2], this allows to
86 have only one linear layer for all the polyseme
87 instead of many ones, each for a polyseme.

88 Image 1, (b) shows this mask process.

89 2.3 One model for coarse and fine

90 An aspect of the given task that has been
91 assumed, is the capacity of a model to benefit from
92 the hierarchical organization of the senses.

93 To this purpose have been tested two model's
94 architecture:

- 95 - Fine grained model to deduct coarse grained
96 (*DeductCg – Model 8*)
- 97 - Joint learning of coarse and fine grained and
98 re-deduction of the fine grained (*Fine-*
99 *Coarse – Model 9*)

100 2.3.1 Deduct Cg Architecture

101 The first architecture that allows to deduct
102 automatically the coarse-grained sense, consist in
103 the training of the previous discussed architecture
104 using the fine-grained classifier. Given the
105 prediction of the model for a given word to be
106 disambiguated the mapping file from the dataset is
107 used to get the "super-sense" label: the coarse
108 grained sense.

109

110 2.3.2 Fine-Coarse Architecture

111 Another study has been conducted based on the
112 assumption that both coarse and fine classifications
113 could mutually enhance each other. Unlike the
114 previously discussed approach of ascending the
115 semantic hierarchy, this study takes a different
116 stance. Rather than attempting to move upward, the
117 model is designed to systematically learn the
118 importance of both tasks concurrently over time.

119 In the Experiments section will be discussed in
120 the detail how this idea has been realized.

121 3 Experiments

122 3.1 Dataset

123 Table 1 summarizes dataset statistics: It comprises
124 6 files with sentences in natural language, part-of-

125 speech info, and lemmas. The sense vocabulary
126 uses WordNet and is composed of 2158 CG senses
127 and 4476 FG senses. Additionally, there's a
128 mapping from fine-grained to coarse-grained
129 senses. Training used accuracy, equivalent to
130 micro-F1 for single-sense word classification.

131 3.2 Experiments List and Setup

132 Table 2 reports all the experiment that have been
133 conducted in this study as well as their identifier.

134 The metric that has been used is the F1-score, all
135 the models have been trained for 100 epochs (some
136 experiments for less epoch due to Colab session
137 duration limit but all of them were not improving
138 anymore). The optimizer is Adam. The learning
139 rate has been reduced with a StepLR scheduler
140 with step size = 1 and gamma = 0.1 starting from
141 1e-3.

142 3.3 Model number 1-2

143 These models, trained for CG WSD follow exactly
144 what is explained in the paragraph 2.2 and they
145 both use the implementation 1 for the f .

146 3.4 Model number 3-4-5

147 Also these models are trained for CG WSD but
148 they use instead the implementation 2 for f .

149 3.5 Model number 6

150 In order to leverage the capacity of LSTM of
151 capturing contextual information, in this
152 experiment, a bidirectional LSTM [7] with 2 layers
153 on top of f is applied before to use the L_1 layer:

$$156 \quad logits_{CG} = Mask_{CG} * L_1(LSTM(f))$$

154 Here in this experiment the best basic model (n. 2)
155 of previous experiments is used to generate f .

157 3.6 Model number 7

158 In this experimental setup, employing the identical
159 architecture as model 2, the focus was on training
160 the model to forecast CG senses. This was achieved
161 by enabling finetuning for the final 4 layers of the
162 BERT model. The rationale behind this approach
163 was to investigate the extent to which the
164 specificity of BERT embeddings could influence
165 the model's performance.

166 3.7 Model number 8

167 Building upon the insights provided in paragraph
168 2.3.1, the model undertakes the dual prediction of
169 both Fine-Grained (FG) senses and Coarse-
170 Grained (CG) senses. The calculation of logits,

171 instrumental in loss computation, is executed as
172 follows:

$$\begin{aligned} 173 \quad & \text{logits}_{FG} = \text{Mask}_{FG} * L_2(f) \\ 174 \quad & p_{FG} = \max(\text{softmax}(\text{logits}_{FG})) \\ 175 \quad & p_{CG} = \text{FG2CG}(p_{FG}) \end{aligned}$$

176 FG2CG is the mapping function that use the given
177 map file to convert a FG sense into CG sense.

178 The central premise behind this model is its
179 optimization for enhancing FG accuracy. The
180 rationale underlying this approach is that if the
181 model attains high performance levels in FG tasks,
182 it is possible to excel in the CG context as well. In
183 essence, a proficient model capable of grasping
184 deeper semantic is expected to reach high
185 performance in more abstract semantic task
186 (clearly only when the task is organized in a formal
187 hierarchical way).

188 3.8 Model number 9

189 To encode the idea in paragraph 2.3.2 in a neural
190 architecture, the first step is to consider a Loss
191 function that is the combination of the two losses:

$$192 \quad \text{Loss} = \text{Loss}_{FG} + \text{Loss}_{CG}$$

193 This setup enables the joint optimization of weights
194 both to decrease CG and FG Loss.

195 Regarding the classification layers: just after the
196 BERT encoder as previous discussed, there are two
197 linear layers (L_1 and L_3):

$$\begin{aligned} 198 \quad & \text{logits}_{CG} = \text{Mask}_{CG} * (L_1(f)) \\ 199 \quad & \text{logits}_{FG} = \text{Mask}_{FG} * L_3(\text{logits}_{CG}) \end{aligned}$$

200 The primary classification layer (L_1) initiates the
201 classification process. Subsequently, the L_3 layer
202 comes into play, facilitating the transformation of
203 logits into a mapping of fine-grained senses;

204 Applying the candidate mask to logits_{CG} it
205 possible to get the CG prediction, and in the same
206 way it is possible to get FG prediction applying
207 the other candidate mask to the logits_{FG} . Much
208 like the prior model, the FG prediction is
209 employed to deduce the CG prediction. However,
210 a potential discrepancy between the CG prediction
211 and the deduced CG prediction can arise. In
212 anticipation of this, the deduced coarse-grained
213 predictions are stored, addressing the likelihood of
214 disparities between the two predictions.

215 This model emits as output:

$$\begin{aligned} 216 \quad & p_{FG} = \max(\text{softmax}(\text{logits}_{FG})) \\ 217 \quad & p_{CG} = \max(\text{softmax}(\text{logits}_{FG})) \\ 218 \quad & p_{CG-Ded} = \text{FG2CG}(p_{FG}) \end{aligned}$$

219 3.9 Model number 10

220 Random baseline models which consist of a
221 random choice over the candidate for each word
222 to disambiguate.

223 4 Results

224 The results in Table 3 reveal important insights
225 about the models. Here are some key points:

226 Embedder Impact on F1: The CG F1 scores show
227 moderate variability across models. Surprisingly,
228 the basic version of BERT outperforms the larger
229 version, suggesting that greater complexity
230 doesn't necessarily yield to superior performance.
231 This phenomenon raises question regarding the
232 trade-off between number of parameters,
233 generalization power and robustness of the
234 outcomes.

235 Implementation 2 for f : Although the F1 scores
236 weren't high, using Implementation 2 for f sped up
237 the forward pass due to smaller tensor dimensions.

238 Direct CG Training: Models 2 achieve high CG
239 F1 scores (0.9150), showing that even a basic
240 architecture can handle the task well.

241 Model 8 attains the highest metric value, despite a
242 minimal difference from Model 2. This highlights
243 that focusing on deeper semantic aspects yields
244 the best performance at higher semantic levels.

245 Joint Learning (Model 9): jointly learning
246 hierarchical sense divisions, requires further
247 exploration for comprehensive understanding.
248 This experiment prompts further investigation.

249 In summary, incorporating hierarchical semantic
250 information involves a blend of neural and
251 symbolic approaches, aligning with previous
252 works. Future research will delve into this
253 emerging AI model category.

254

255

256

257

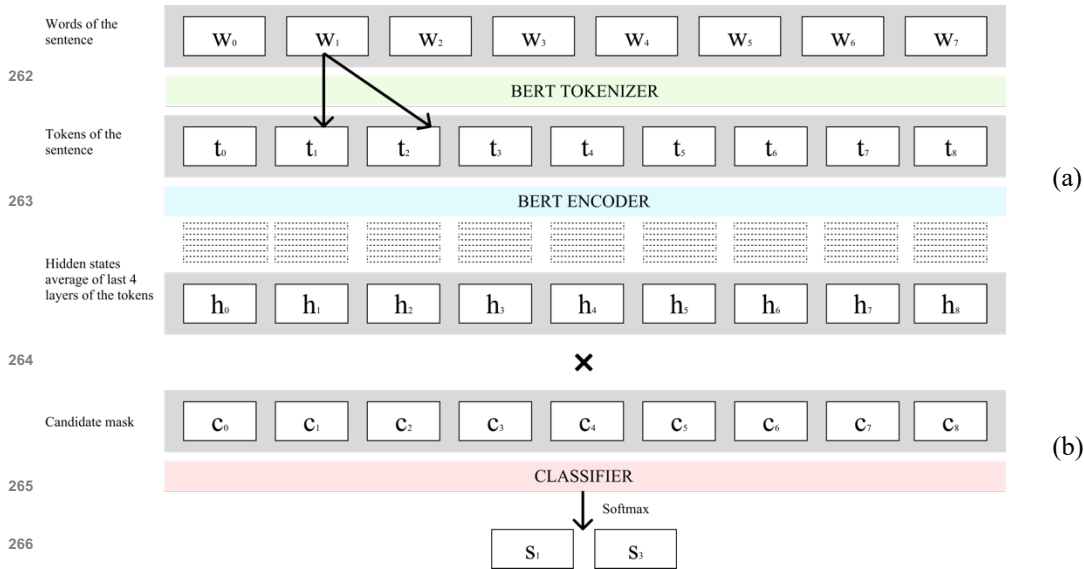
258

259

260

Figure 1: (a) The encoding part of the model. (b) The mask process.

261



266

	Training	Validation	Test
# Sentences	12339	685	686
# sentence with 1 candidate	7723	451	425
Sentence length mean	40.80	40.95	41.22

Table 1: Statistic for the training validation and test data plus the count of the number of words with 1 candidate and the mean of words length.

ID	CG F1	FG F1
1	0.9094	/
2	0.9150	/
3	0.8752	/
4	0.8629	/
5	0.8668	/
6	0.8920	/
7	0.9045	/
8	0.9151	0.8041
9	0.9143	0.7976
10	0.68	0.52

Table 3: Metric value for the models discussed in the experiments section. The ID column indicated the model identifier. For model 8, is reported the deducted CG because is always higher than the predicted CG metric. Model 10 represent the random Baseline model

ID	Granularity	Embedder
1	CG	Bert Large Cased
2	CG	Bert Base Cased
3	CG	Bert Large Cased
4	CG	Bert Large uncased
5	CG	Bert Base Cased
6	CG	Bert Base Cased + LSTM
7	CG	Bert Based Cased
8	FG&CG	Bert Based Cased
9	FG&CG	Bert Based Cased

Table 2: Experiment List: 1,2,6,7 with implementation 1 of the embedding, 2-3-4 with implementation 2, 8 with implementation 3 and 9 with implementation 4; Exp 8 has the last 4 layers of Bert finetuned.

267 References

- 268 1. “Word Sense Disambiguation with Recurrent
269 Neural Networks” – Alexander Popov,
270 Linguistic Modelling Department IICT-BAS
- 271 2. “Using BERT for Word Sense
272 Disambiguation” - Jiaju Du, Fanchao Qi,
273 Maosong Sun
- 274 3. “Automatic sense disambiguation using
275 machine readable dictionaries: how to tell a
276 pine cone from an ice cream cone”- Michael
277 Lesk. 1986.
- 278 4. “An Enhanced Lesk Word Sense
279 Disambiguation Algorithm through a
280 Distributional Semantic Model” - Basile,
281 Caputo, Semeraro 2014
- 282 5. [BERT: Pre-training of Deep Bidirectional
283 Transformers for Language Understanding]
284 (Devlin et al., NAACL 2019)
- 285 6. [Breaking Through the 80% Glass Ceiling:
286 Raising the State of the Art in Word Sense
287 Disambiguation by Incorporating Knowledge
288 Graph Information](Bevilacqua & Navigli,
289 ACL 2020)
- 290 7. Hochreiter, Sepp & Schmidhuber, Jürgen.
291 (1997). Long Short-term Memory. Neural
292 computation. 9. 1735-80.
293 10.1162/neco.1997.9.8.1735.

294