

Context	Objectives
<ul style="list-style-type: none"> <li>Extraction of lexical semantic relations in English               <ul style="list-style-type: none"> <li>Paradigmatic relations and more particularly synonyms</li> </ul> </li> <li>General method: distributional thesauri               <ul style="list-style-type: none"> <li>For each target word: retrieval of its first neighbors according to the distributional representation of words</li> </ul> </li> <li>According to the principles of distributional analysis (Harris, 1954)               <ul style="list-style-type: none"> <li>first neighbors = synonyms of the target words</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Adaptation to contextual embeddings (ELMo, BERT...)               <ul style="list-style-type: none"> <li>Static embeddings = word-level representations</li> <li>Contextual embeddings = token-level representations</li> </ul> </li> </ul> <p style="text-align: center;">➔ Building of word representations from contextual language models</p>

Existing methods	Proposed method
<p><b>Averaging of token embeddings (Bommasani et al., 2020)</b></p> <ul style="list-style-type: none"> <li>Encoding of sentences containing the target word with a contextual model</li> <li>Word representation = average of embeddings extracted for the target word in all sentences</li> </ul> <p><b>Principal Component Analysis (Ethayarajh, 2019)</b></p> <ul style="list-style-type: none"> <li>Word representation = first principal component of token representations</li> </ul> <p><b>Token without context (Vulić et al., 2020)</b></p> <ul style="list-style-type: none"> <li>Encoding a pseudo-sentence made of only one word</li> </ul>	<p><b>Starting point</b></p> <ul style="list-style-type: none"> <li>For each target word: set of token representations = in-context word representations</li> </ul> <p><b>General idea</b></p> <ul style="list-style-type: none"> <li>Exploiting the diversity of token representations for building a more general representation at the word level</li> </ul> <p><b>3 main steps</b></p> <ul style="list-style-type: none"> <li>Selection of a focused set of token representations for each target word               <ul style="list-style-type: none"> <li>for limiting the heterogeneity of representations</li> </ul> </li> <li>Generalization of selected token representations</li> <li>Aggregation of generalized token representations</li> </ul>

Representation selection	Representation generalization and aggregation
<p><b>Starting point</b></p> <ul style="list-style-type: none"> <li>Random selection of <math>N</math> sentences from the AQUAINT corpus (news articles) for each target word</li> <li>Target words = nouns</li> <li><math>N</math> between 10 and 250 according to the frequency of the target nouns</li> </ul> <p><b>Hypothesis</b></p> <ul style="list-style-type: none"> <li>(McCarthy et al., 2004): dominant sense for each word in a corpus</li> <li>Average token representation = representation of the dominant sense of the associated word</li> </ul> <p><b>4 selection strategies</b></p> <ul style="list-style-type: none"> <li>Random</li> <li>Closest to the dominant sense</li> <li>Farthest to the dominant sense</li> <li>Uniform distribution in terms of proximity with the dominant sense</li> </ul> <p style="text-align: center;">→ 10 token representations / word</p>	<p><b>Generalization principles</b></p> <ul style="list-style-type: none"> <li>Adaptation of the token representations of a word to bring them closer together</li> <li>Similar to knowledge injection methods such as <i>retrofitting</i> (Faruqui et al., 2015)               <ul style="list-style-type: none"> <li>as if the token representations of a word were linked by similarity relations</li> </ul> </li> </ul> <p><b>2 main variants</b></p> <ul style="list-style-type: none"> <li><b>gen+agg</b>: generalization, then aggregation               <ul style="list-style-type: none"> <li>generalization                   <ul style="list-style-type: none"> <li>pseudo-similarity relation between each pair of token representations</li> <li>application of a retrofitting-like algorithm: PARAGRAM (Wieting et al., 2015)</li> </ul> </li> <li>aggregation: averaging of generalized token representations</li> </ul> </li> <li><b>agg+gen</b>: aggregation, then generalization               <ul style="list-style-type: none"> <li>aggregation: averaging of token representations</li> <li>generation of pseudo-similarity relations                   <ul style="list-style-type: none"> <li><b>agg+gen<sub>all</sub></b>: between each pair of token representations, including the aggregate</li> <li><b>agg+gen<sub>agg</sub></b>: only between the aggregate and the initial token representations</li> </ul> </li> <li>application of the retrofitting-like algorithm to all representations</li> </ul> </li> </ul> <p>➔ final word representation = aggregate</p> <div style="text-align: right;"> <p>— pseudo-similarity relation</p> <p>● token representation</p> <p>● aggregate</p> </div>

Experiments and evaluation	Evaluation of the proposed method	Table 1																																																												
<p><b>Large-scale intrinsic evaluation</b></p> <ul style="list-style-type: none"> <li>Gold Standard: WordNet's synonyms</li> <li>Evaluated target words = 10,305 nouns               <ul style="list-style-type: none"> <li>frequency: large spectrum</li> </ul> </li> <li>Information Retrieval (IR) paradigm               <ul style="list-style-type: none"> <li>target word <math>\equiv</math> query; neighbors <math>\equiv</math> docs</li> <li>IR measures: MAP, R-precision, precision@{1,2,5}</li> </ul> </li> <li>Evaluated models               <ul style="list-style-type: none"> <li>BERT uncased: token representation = average representation of its WordPieces</li> <li>CharacterBERT (CBERT): influence of WordPieces</li> <li>fastText: static embeddings (Skip-gram model) from Wikipedia</li> </ul> </li> </ul>	<p><b>Evaluation of the proposed method</b></p> <ul style="list-style-type: none"> <li>Reference methods: *-avg, *-pca, *-iso               <ul style="list-style-type: none"> <li>*-iso &lt; *-avg = *-pca ~ fastText</li> </ul> </li> <li>Proposed method &gt; reference methods and static embeddings               <ul style="list-style-type: none"> <li>for both BERT and CharacterBERT</li> </ul> </li> <li>aggregation+generalization &gt; generalization+aggregation</li> <li>agg+gen<sub>agg</sub> ~ agg+gen<sub>all</sub> <ul style="list-style-type: none"> <li>preference for agg+gen<sub>agg</sub>: requires fewer computations</li> </ul> </li> <li>Similar improvement for CharacterBERT and BERT</li> </ul>	<table border="1"> <thead> <tr> <th></th> <th>R<sub>prec</sub></th> <th>MAP</th> <th>P@1</th> <th>P@2</th> <th>P@5</th> </tr> </thead> <tbody> <tr><td>CBERT-avg</td><td>15.6</td><td>18.0</td><td>22.0</td><td>15.9</td><td>9.2</td></tr> <tr><td>CBERT-pca</td><td>15.6</td><td>17.9</td><td>22.0</td><td>15.9</td><td>9.2</td></tr> <tr><td>CBERT-gen+agg</td><td>16.1</td><td>18.6</td><td>22.6</td><td>16.3</td><td>9.5</td></tr> <tr><td>CBERT-agg+gen<sub>all</sub></td><td><b>16.3</b></td><td><b>18.9</b></td><td><b>22.8</b></td><td><b>16.5</b></td><td>9.7</td></tr> <tr><td>CBERT-agg+gen<sub>agg</sub></td><td><b>16.3</b></td><td>18.8</td><td><b>22.8</b></td><td>16.4</td><td>9.6</td></tr> <tr><td>BERT-avg</td><td>15.6</td><td>17.9</td><td>21.8</td><td>16.0</td><td>9.5</td></tr> <tr><td>BERT-iso (L0)</td><td>14.0</td><td>15.8</td><td>19.2</td><td>14.6</td><td>8.7</td></tr> <tr><td>BERT-agg+gen<sub>agg</sub></td><td>16.2</td><td>18.8</td><td>22.5</td><td>16.1</td><td><b>10.1</b></td></tr> <tr><td>fastText</td><td>15.5</td><td>18.4</td><td>21.9</td><td>15.7</td><td>9.2</td></tr> </tbody> </table> <p>Best layer for CharacterBERT = L4    Random selection for token representations Best layer for BERT = L5</p>		R <sub>prec</sub>	MAP	P@1	P@2	P@5	CBERT-avg	15.6	18.0	22.0	15.9	9.2	CBERT-pca	15.6	17.9	22.0	15.9	9.2	CBERT-gen+agg	16.1	18.6	22.6	16.3	9.5	CBERT-agg+gen <sub>all</sub>	<b>16.3</b>	<b>18.9</b>	<b>22.8</b>	<b>16.5</b>	9.7	CBERT-agg+gen <sub>agg</sub>	<b>16.3</b>	18.8	<b>22.8</b>	16.4	9.6	BERT-avg	15.6	17.9	21.8	16.0	9.5	BERT-iso (L0)	14.0	15.8	19.2	14.6	8.7	BERT-agg+gen <sub>agg</sub>	16.2	18.8	22.5	16.1	<b>10.1</b>	fastText	15.5	18.4	21.9	15.7	9.2
	R <sub>prec</sub>	MAP	P@1	P@2	P@5																																																									
CBERT-avg	15.6	18.0	22.0	15.9	9.2																																																									
CBERT-pca	15.6	17.9	22.0	15.9	9.2																																																									
CBERT-gen+agg	16.1	18.6	22.6	16.3	9.5																																																									
CBERT-agg+gen <sub>all</sub>	<b>16.3</b>	<b>18.9</b>	<b>22.8</b>	<b>16.5</b>	9.7																																																									
CBERT-agg+gen <sub>agg</sub>	<b>16.3</b>	18.8	<b>22.8</b>	16.4	9.6																																																									
BERT-avg	15.6	17.9	21.8	16.0	9.5																																																									
BERT-iso (L0)	14.0	15.8	19.2	14.6	8.7																																																									
BERT-agg+gen <sub>agg</sub>	16.2	18.8	22.5	16.1	<b>10.1</b>																																																									
fastText	15.5	18.4	21.9	15.7	9.2																																																									

Evaluation of selection methods for token representations	Word frequency analysis	Table 3																																																																		
<ul style="list-style-type: none"> <li>Results for agg+gen<sub>agg</sub></li> <li>Equivalence of all strategies</li> <li>Limited effect of the contextualization of embeddings</li> </ul> <p>➔ uniform: « logical » choice as a deterministic strategy</p> <table border="1"> <thead> <tr> <th></th> <th>R<sub>prec</sub></th> <th>MAP</th> <th>P@1</th> <th>P@2</th> <th>P@5</th> </tr> </thead> <tbody> <tr><td>random</td><td>16.3</td><td>18.8</td><td>22.8</td><td>16.4</td><td>9.6</td></tr> <tr><td>uniform</td><td>16.4</td><td><b>18.9</b></td><td><b>22.9</b></td><td><b>16.7</b></td><td>9.7</td></tr> <tr><td>farthest</td><td><b>16.5</b></td><td><b>18.9</b></td><td><b>22.9</b></td><td>16.6</td><td><b>9.8</b></td></tr> <tr><td>closest</td><td>16.3</td><td>18.8</td><td><b>22.9</b></td><td>16.4</td><td>9.6</td></tr> </tbody> </table> <p style="text-align: center;">Table 2</p>		R <sub>prec</sub>	MAP	P@1	P@2	P@5	random	16.3	18.8	22.8	16.4	9.6	uniform	16.4	<b>18.9</b>	<b>22.9</b>	<b>16.7</b>	9.7	farthest	<b>16.5</b>	<b>18.9</b>	<b>22.9</b>	16.6	<b>9.8</b>	closest	16.3	18.8	<b>22.9</b>	16.4	9.6	<ul style="list-style-type: none"> <li>Split of results according to the median frequency of target words</li> <li>High-frequency words               <ul style="list-style-type: none"> <li>contextual &gt; static</li> <li>impact of contextualization on polysemy</li> </ul> </li> <li>Low-frequency words: opposite trend</li> <li>Fusion of thesauri (CombSum strategy): significant improvement / static and contextual</li> </ul>	<table border="1"> <thead> <tr> <th></th> <th>R<sub>prec</sub></th> <th>MAP</th> <th>P@1</th> <th>P@2</th> <th>P@5</th> </tr> </thead> <tbody> <tr><td>uniform<sub>high</sub></td><td>18.0</td><td>20.5</td><td>26.6</td><td>19.8</td><td>12.0</td></tr> <tr><td>uniform<sub>low</sub></td><td>14.8</td><td>17.3</td><td>19.3</td><td>13.4</td><td>7.3</td></tr> <tr><td>fastText<sub>high</sub></td><td>14.5</td><td>16.8</td><td>22.0</td><td>15.9</td><td>9.7</td></tr> <tr><td>fastText<sub>low</sub></td><td>16.4</td><td>20.0</td><td>21.8</td><td>15.4</td><td>8.6</td></tr> <tr><td>CombSum</td><td>18.6</td><td>21.5</td><td>25.9</td><td>19.0</td><td>11.1</td></tr> </tbody> </table> <p style="text-align: center;">Table 3</p>		R <sub>prec</sub>	MAP	P@1	P@2	P@5	uniform <sub>high</sub>	18.0	20.5	26.6	19.8	12.0	uniform <sub>low</sub>	14.8	17.3	19.3	13.4	7.3	fastText <sub>high</sub>	14.5	16.8	22.0	15.9	9.7	fastText <sub>low</sub>	16.4	20.0	21.8	15.4	8.6	CombSum	18.6	21.5	25.9	19.0	11.1
	R <sub>prec</sub>	MAP	P@1	P@2	P@5																																																															
random	16.3	18.8	22.8	16.4	9.6																																																															
uniform	16.4	<b>18.9</b>	<b>22.9</b>	<b>16.7</b>	9.7																																																															
farthest	<b>16.5</b>	<b>18.9</b>	<b>22.9</b>	16.6	<b>9.8</b>																																																															
closest	16.3	18.8	<b>22.9</b>	16.4	9.6																																																															
	R <sub>prec</sub>	MAP	P@1	P@2	P@5																																																															
uniform <sub>high</sub>	18.0	20.5	26.6	19.8	12.0																																																															
uniform <sub>low</sub>	14.8	17.3	19.3	13.4	7.3																																																															
fastText <sub>high</sub>	14.5	16.8	22.0	15.9	9.7																																																															
fastText <sub>low</sub>	16.4	20.0	21.8	15.4	8.6																																																															
CombSum	18.6	21.5	25.9	19.0	11.1																																																															