

## A large reproducible benchmark of ontology-based methods and word embeddings for word similarity

Juan J. Lastra-Díaz <sup>1a</sup>, Josu Goikoetxea <sup>2b</sup>, Mohamed Ali Hadj Taieb <sup>3c</sup>, Ana Garcia-Serrano <sup>4a</sup>, Mohamed Ben Aouicha <sup>5c</sup>, Eneko Agirre <sup>6b</sup>

<sup>a</sup>*NLP & IR Research Group, ETSI de Informática (UNED)  
Universidad Nacional de Educación a Distancia, Juan del Rosal 16, 28040 Madrid (Spain)*

<sup>b</sup>*IXA NLP group, Faculty of Informatics, UPV/EHU  
Manuel Lardizabal 1 (20018), Donostia, Basque Country*

<sup>c</sup>*Faculty of Sciences of Sfax, Tunisia*

---

### Abstract

This work is a companion reproducibility paper of the experiments and results introduced by Lastra-Díaz et al. [44], which is based on the evaluation of a companion reproducibility dataset with the HESML V1R4 library and the long-term reproducibility tool called Rezip. Human similarity and relatedness judgements between concepts underlie most of cognitive capabilities, such as categorisation, memory, decision-making and reasoning. For this reason, the research on methods for the estimation of the degree of similarity and relatedness between words and concepts has received a lot of attention in the fields of artificial intelligence and cognitive sciences. However, despite the huge research effort done, there is a lack of a self-contained, reproducible and extensible collection of benchmarks which being amenable to become a de facto standard for large scale experimentation in this line of research. In order to bridge this reproducibility gap, this work introduces a set of reproducible experiments on word similarity and relatedness by providing a detailed reproducibility protocol together with a set of software tools and a self-contained reproducibility dataset which allow that all experiments and results in our aforementioned work to be reproduced exactly. Our aforementioned primary work introduces the largest, detailed and reproducible experimental survey on word similarity and relatedness reported in the literature, which is based on the implementation of all evaluated methods into a same software platform. Our reproducible experiments evaluate most of methods in the families of ontology-based semantic similarity measures and word embedding models. Finally, we detail how to extend our experiments to evaluate other unconsidered experimental setups.

### Keywords:

Ontology-based semantic similarity measures, Word embeddings, Information Content models, Reproducible benchmark, HESML, Rezip

---

<sup>1</sup>jlustra@invi.uned.es (corresponding author)

<sup>2</sup>josu.goikoetxea@ehu.eus

<sup>3</sup>mohamedali.hadjtaieb@gmail.com

<sup>4</sup>agarcia@lsi.uned.es

<sup>5</sup>mohamed.benaouicha@fss.usf.tn

<sup>6</sup>e.agirre@ehu.eus

---

## 1. Introduction

Human similarity and relatedness judgements between concepts underlie most of cognitive capabilities, such as categorisation, memory, decision-making and reasoning. Thus, the proposal of methods for the estimation of the degree of similarity and relatedness between words and concepts has been a very active line of research in the fields of Artificial Intelligence (AI), Natural Language Processing (NLP) and Information Retrieval (IR). For this reason, authors has been largely involved in this line of research during the last decade. For instance, we have proposed new ontology-based semantic similarity measures [36, 39, 33, 34, 25, 24, 23, 22, 4, 3, 6], Information Content (IC) models [35, 37, 27, 26, 6], word embeddings [18, 20, 19, 21], distributional semantics measures [68, 1, 7], semantic measure libraries [43, 5], reproducibility resources [40, 45], word similarity benchmarks [1], reproducible experiments on word similarity based on WordNet [34, 41] and reproducible benchmarks between semantic measures libraries [34, 38].

Main approaches on word similarity and relatedness proposed in the literature can be categorised in two large families as follows: (1) Ontology-based semantic similarity Measures (OM), and (2) distributional measures whose most recent and successful methods are based on Word Embedding (WE) models. Precisely, our primary work [44] introduces a deep experimental study on both aforementioned families of methods encouraged by the recent advances in the family of WE models, which is based on the implementation and evaluation of all methods in a same software platform based on HESML V1R4 [42] and WordNet 3.0 [50], as well as their subsequent recording with the Rezip long-term reproducibility tool [13]. Before the publication of this work, the only large reproducible experimental surveys on word similarity reported in the literature were those introduced by Lastra-Díaz and García-Serrano [36, 35, 37] in another reproducibility paper [43] belonging to this same reproducibility initiative [12], in which we also find other works such as those introduced by Wolke et al. [66] and Fariña et al. [14]. However, there is no neither joint reproducible benchmarks on word embeddings and ontology-based semantic similarity measures nor other ones evaluating the latest family of methods on so large count of datasets as those evaluated by our primary work [44].

Reproducibility of methods and research results in the field of NLP has become a serious problem which severely hampers any research effort and the smooth integration of newcomers in the field, which is especially hard for most graduate students who start their scientific career in this aforementioned field. This reproducibility gap was already highlighted in a pioneering work by Pedersen [54], being subsequently confirmed by Fokkens et al. [16] by evaluating several works in the same line of research tackled herein. More recently, Branco et al. [9] introduce a call for reproducibility submissions in a known NLP journal to bridge the aforementioned reproducibility gap. We subscribed to this reproducibility alarm by adopting as basic norm the detailed replication of all methods evaluated in our papers, as well as the warning on many contradictory or unreproducible results in a series of papers in this line of research, such as the works introduced by Lastra-Díaz and García-Serrano [36, 35, 37], Lastra-Díaz et al. [43] and our primary paper [44]. In a recent and valuable reproducibility

48 study in the field of NLP, Wieling et al. [64, p.641] found that only a third part  
49 of the published works (36.2%) provided their source code; however, they found  
50 by evaluating a random sample of ten works that only a tenth part of the former  
51 group could be reproduced exactly. Thus, this later finding yields an alarming  
52 ratio of only a 3.62% of reproducible works in this aforementioned study. For  
53 all reasons above, we subscribe both the reproducible manifesto [53] for a repro-  
54 ducible science and reproducibility initiative lead by Information Systems [12],  
55 as well as the slow science manifesto<sup>7</sup> for a reflective research. Finally, we make  
56 our own the words of Pedersen [54, p.470]: “we might one day only accept for  
57 publication articles that are accompanied by working software that allows for  
58 immediate and reliable reproduction of results”.

59 The aim of this work is to introduce a detailed experimental setup based on  
60 a collection of publicly available software tools [42] and reproducibility resources  
61 [45, 46], which are provided as supplementary material, with the aim of exactly  
62 reproducing all experiments and results reported in our primary work [44].

### 63 *1.1. Main motivation*

64 Our first motivation is to introduce a self-contained and easily reproducible  
65 set of experiments on word similarity and relatedness which allow to reproduce  
66 all experiments, results and conclusions introduced by our primary work [44]  
67 exactly.

68 A second motivation is the lack of a self-contained, reproducible and exten-  
69 sible collection of benchmarks on word similarity and relatedness which jointly  
70 evaluate the most recent methods on the families of ontology-based semantic  
71 similarity measures and word embedding models on a same software platform,  
72 and consequently, being amenable to become a de facto standard for large scale  
73 experimentation in this line of research. Despite the huge research effort done  
74 during the last decades, such as witnessed by the plethora of methods reviewed  
75 and evaluated in our primary work [44], there is still a lack of a fully auto-  
76 matic, reproducible and extensible collection of benchmarks which make the  
77 evaluation and development of word similarity and relatedness methods easier.  
78 In general, there is a lack of reproducibility resources in this line of research  
79 which was partially bridged by the introduction of several semantic measures  
80 libraries, such as SML [29], SISR [5] and the most recent HESML [43] which is  
81 the largest and efficient among them, in addition to provide self-contained and  
82 easily reproducible experiments by the first time. Likewise, the reproducible  
83 experiments and reproducibility datasets introduced by Lastra-Díaz et al. [43]  
84 and Lastra-Díaz and García-Serrano [40, 41, 38] respectively have allowed by  
85 the first time to reproduce a set of large experimental surveys on ontology-based  
86 semantic similarity measures based on WordNet [36, 35, 37] exactly. However,  
87 recent and fast advances in the family of word embedding models together with  
88 the active research on ontology-based methods have raised the need to carry-out  
89 joint evaluations of both families of methods in a large set of benchmarks to  
90 elucidate the state of the problem.

---

<sup>7</sup><http://slow-science.org/>

91 *1.2. Definition of the problem and contributions*

92 This work tackles the problem of designing and implementing for the first  
93 time a self-contained and easily reproducible set of experiments on word similar-  
94 ity and relatedness of the families of ontology-based semantic similarity measures  
95 based on WordNet and word embeddings by providing a very detailed repro-  
96 ducibility protocol together with a set of software tools [42] and a companion  
97 reproducibility dataset [45] which is publicly available at [46]. In addition, we  
98 detail how our reproducible experiments could be extended for setting up and  
99 evaluating unconsidered experimental setups including other datasets, word em-  
100 beddings or ontology-based semantic similarity measures.

101 The rest of the paper is structured as follows. Section 2 introduces HESML  
102 library [43] and Rezip tool [12] which set the software platform originally  
103 used to run all experiments introduced herein and our long-term reproducibility  
104 platform respectively. Section 3 introduces the new reproducible experiments  
105 on word similarity whilst section 4 details how them can be extended or created  
106 new ones from scratch. Finally, we introduce our conclusions and future work.

107 **2. Background on HESML and Rezip**

108 *HESML* [43] is a self-contained Java software library of semantic measures  
109 based on WordNet whose latest version, called HESML V1R4 [42], also supports  
110 the evaluation of pre-trained word embedding models, such as those introduced  
111 by Mikolov et al. [49], Pennington et al. [57], Schwartz et al. [61], Wieting  
112 et al. [65], Goikoetxea et al. [20], Bojanowski et al. [8], Agirre and Soroa [2],  
113 Camacho-Collados et al. [11] and Mrkšić et al. [52]. HESML is a self-contained  
114 experimentation platform on word similarity and relatedness which is especially  
115 well suited to run large experimental surveys by supporting the execution of au-  
116 tomatic reproducible experiment files on word similarity based on a XML-based  
117 file format (\*.exp). Despite the latest version of HESML only supports Word-  
118 Net, it could be easily extended to manage other ontologies by implementing  
119 the proper parsers as detailed by Lastra-Díaz et al. [43]. HESML library has  
120 been completely developed in NetBeans 8 and Java 8, being distributed with  
121 three WordNet versions whilst *HESMLclient* is a complementary Java console  
122 program whose aim is to run word similarity experiments by calling HESML  
123 functionality. For a detailed introduction to HESML, we refer the reader to  
124 its introductory paper [43]. Table 1 shows a summary of technical and legal  
125 information of the latest HESML version used in our experiments.

126 On the other hand, Rezip is a virtualization tool introduced by Chirigati  
127 et al. [13], whose aim is to warrant the exact replication of experimental results  
128 onto different systems from that originally used in their creation. Rezip  
129 captures all the program dependencies and is able to reproduce the packaged  
130 experiments on any host platform, regardless of the hardware and software con-  
131 figuration used in their creation. Thus, Rezip warrants the reproduction  
132 of the experiments introduced herein in the long term. Other valuable feature  
133 of Rezip is that it allows to modify the input files of any Rezip package  
134 with the aim of evaluating a set of experiments using originally unconsidered  
135 methods, configuration parameters or datasets. Rezip supports main virtu-  
136 alization platforms as Docker and VirtualBox; however, our preferred option is  
137 Docker. For a comparison of these two types of virtualization platforms, we

HESML software library	Description
Current code version.	V1R4
Legal Code License.	Creative Commons By-NC-SA 4.0
Permanent code repository used for this version.	<a href="http://dx.doi.org/10.17632/t87s78dg78.4">http://dx.doi.org/10.17632/t87s78dg78.4</a>
GitHub repository	<a href="https://github.com/jjlastra/HESML.git">https://github.com/jjlastra/HESML.git</a>
Software code languages and tools.	Java 8, Java SE DevKit 8, NetBeans 8.0 or higher
Compilation requirements and operating systems.	Java SE Dev Kit 8, NetBeans 8.0 or higher and any Java-compliant operating system.
Documentation and source code examples	Sample source code in the HESMLclient program.
Community forum for questions.	hesml+subscribe@googlegroups.com, hesml+unsubscribe@googlegroups.com

Table 1: Technical and legal information of the latest version of the HESML software library [43] used in our experiments.

Reprozip tool	Description
Current version	1.0.16
Web site	<a href="https://www.reprozip.org">https://www.reprozip.org</a>
Supported platforms	Linux, Windows and MacOS

Table 2: Technical and access information of Reprozip long-term reproducibility tool [12].

138 refer the reader to the survey introduced by Merkel [48], in which the author  
139 introduces Docker and compares it with classic Virtual Machines (VM) such  
140 as VirtualBox. Finally, Reprozip also simplifies the generation, packaging and  
141 execution of Docker-based experiments. For all reasons above, we encourage  
142 the research community to use Reprozip as a long-term reproducibility backup.  
143 Table 2 shows a summary of technical and access information of the Reprozip  
144 reproducibility tool.

### 145 3. The Reproducible Experiments on word similarity

146 The aim of this section is to introduce a set of detailed experimental setups  
147 in order to replicate the methods and experiments introduced by our primary  
148 work [44] exactly. Section 3.1 details the experimental setup for the implemen-  
149 tation of our experiments in our primary work [44], then section 3.2 details the  
150 minimal system requirements for the testing platforms with the aim of running  
151 our reproducible experiments. Likewise, section 3.2 reports the running times  
152 obtained by the authors and reviewers in the evaluation of our reproducible  
153 experiments in different testing platforms. Section 3.3 details the procedure  
154 for obtaining and compiling HESML source code, as well as running its pre-  
155 compiled jar files. We note that it is unneeded to compile the HESML source  
156 code to run the experiments because the HESML distribution already includes  
157 pre-compiled versions of the HESMLclient program with the latest HESML ver-  
158 sion. Next, section 3.4 introduces the method to run our experiments which are  
159 based on the running of HESMLclient program whilst section 3.5 introduces

160 our long-term reproducibility method based on Rezip. Finally, section 3.6  
 161 introduces the automated data analysis carried-out to process the raw similar-  
 162 ity values generated by our experiments and computing all evaluation metrics  
 163 reported in our aforementioned primary work [44], as well as a report in HTML  
 164 file format showing all data tables generated from our raw data.

### 165 3.1. Experimental setup in our primary paper

166 All experiments carried-out in our primary paper [44] were implemented in  
 167 *HESML V1R4* [42] by running *HESMLclient* program with a reproducible ex-  
 168 periment file in XML-based (\*.exp) file format which encodes the evaluation of  
 169 all semantic measures in all datasets as listed in table 4. Experimental setup  
 170 and software platform used to implement all our experiments is detailed in [44,  
 171 figure 3]. HESML V1R4 implements all ontology-based semantic similarity mea-  
 172 sures based on WordNet 3.0 as well as the evaluation of all pre-trained word  
 173 embedding models evaluated in our benchmarks. In addition, execution of our  
 174 experiments was recorded into a long-term reproducibility Rezip file, called  
 175 “*WN\_ontology\_measures\_vs\_embeddings.rpz*”, which is part of our companion re-  
 176 producibility dataset [45], being publicly available at [46]. Our aforementioned  
 177 Rezip file can be reproduced in any Rezip compliant platform<sup>8</sup> as detailed  
 178 in section 3.5. Thus, all our methods, experiments and results can be repro-  
 179 duced using two different software platforms and methods as detailed in table 3.  
 180 First reproducibility method is based on the execution of *HESMLclient* program  
 181 whilst second one is based on the execution of the aforementioned Rezip file.

Software used	Supported reproducibility methods
HESMLclient [42]	You should download HESML V1R4 [42] and a supplementary ZIP file containing the collection of pre-trained word embedding files ( <i>WordEmbeddings.zip</i> [46]), and then running <i>HESMLclient</i> with the reproducible file as input, as detailed in section 3.4.
ReproUnzip [13]	You should download our supplementary Rezip file [46] and setting up and running Rezip as detailed in section 3.5.

Table 3: Our two methods to reproduce all experiments and results introduced by our primary work [44]. HESMLclient method is that originally used to run our experiments in our primary work, whilst ReproUnzip provides a long-term reproducibility method regardless the original testing platform used to run our experiments.

182 The two reproducibility methods cited in table 3 were introduced in the  
 183 introductory HESML paper [43] which introduces a detailed protocol to repro-  
 184 duce all experiments, results, data tables and figures reported in three papers  
 185 previously introduced by Lastra-Díaz and García-Serrano [35, 36, 37], as well  
 186 as the benchmarks between semantic measures libraries reported in [43]. All  
 187 experiments detailed herein were originally implemented on an UBUNTU 16.04

<sup>8</sup><https://www.reprozip.org/>

188 virtual computer with 8 Gb of RAM and 100 Gb of disk space called UBUNTU-  
 189 base1 as detailed in table 5. However, it could be reproduced in any Java 8 or  
 190 Reprounzip compliant platform by using any of the two aforementioned meth-  
 191 ods above, which includes most Linux-based, MacOS-based and Windows-based  
 192 platforms. For this reason, our experiments have been successfully reproduced  
 193 using both HESMLclient and Reprounzip methods (see figure 1) in all testing  
 194 platforms detailed in table 5 with the running times reported in table 6.

195 Figure 1 shows the three reproducibility workflows introduced herein which  
 196 are defined by the selection of one of the two reproducibility methods shown in  
 197 table 3 with a specific testing platform. HESML distribution includes the pre-  
 198 compiled version of HESML V1R4 and HESMLclient .jar files, thus any reader  
 199 interested in reproducing our experiments can directly follow the setup instruc-  
 200 tions in tables 7 and 8, and subsequently running the experiments as detailed  
 201 in tables 10 and 11. On the other hand, table 4 shows the full collection of re-  
 202 producible experiments encoded by the “*benchmark\_survey.exp*” file (see figure  
 203 2), as well as the corresponding raw output files that are generated during its  
 204 execution whose subsequent processing allows to reproduce the results reported  
 205 in our primary paper [44] exactly, as detailed in section 3.6.

Word similarity and relatedness benchmarks reproduced herein [44]

Dataset	Raw output files generated by our experiments
MC28 [51]	raw_similarity_values_MC28_dataset.csv
RG65 [60]	raw_similarity_values_RG65_dataset.csv
PS <sub>full</sub> [58]	raw_similarity_values_PSfull_dataset.csv
Agirre201 [1]	raw_similarity_values_Agirre201_lowercase_dataset.csv
SimLex665 [30]	raw_similarity_values_SimLex665_dataset.csv
MTurk771 [28]	raw_similarity_values_MTurk771_dataset.csv
MTurk287/235 [59]	raw_similarity_values_MTurk287-235_dataset.csv
WS353Rel [15]	raw_similarity_values_WS353Rel_dataset.csv
Rel122 [63]	raw_similarity_values_Rel122_dataset.csv
WS353Full [15]	raw_similarity_values_WS353Full_dataset.csv
SimLex111 [30]	raw_similarity_values_SimLex111_dataset.csv
SimLex222 [30]	raw_similarity_values_SimLex222_dataset.csv
SimLex999 [30]	raw_similarity_values_SimLex999_dataset.csv
SimVerb3500 [17]	raw_similarity_values_SimVerb3500_dataset.csv
MEN [10]	raw_similarity_values_MEN_dataset.csv
YP130 [67]	raw_similarity_values_YP130_dataset.csv
RW2034 [47]	raw_similarity_values_RareWords2034_dataset.csv
RW1401 [47]	raw_similarity_values_RareWords1401_dataset.csv
SCWS [32]	raw_similarity_values_SCWS1994_dataset.csv

Table 4: Collection of raw output files generated by the execution of the “*benchmark\_survey.exp*” reproducible experiment file by any of the two aforementioned reproducibility methods. Each raw output file contains the raw similarity or relatedness values returned for each word pair by each semantic measure. These raw output files are subsequently processed by a R-language script to produce the final data tables shown in our primary paper [44], as detailed in section 3.6. For further details on the datasets above, we refer the reader to our primary paper [44, table 3].

206 HESML V1R4 distribution [42] contains all source files and pre-compiled  
 207 versions of the *HESML-V1R4.jar* library and *HESMLclient.jar* Java console

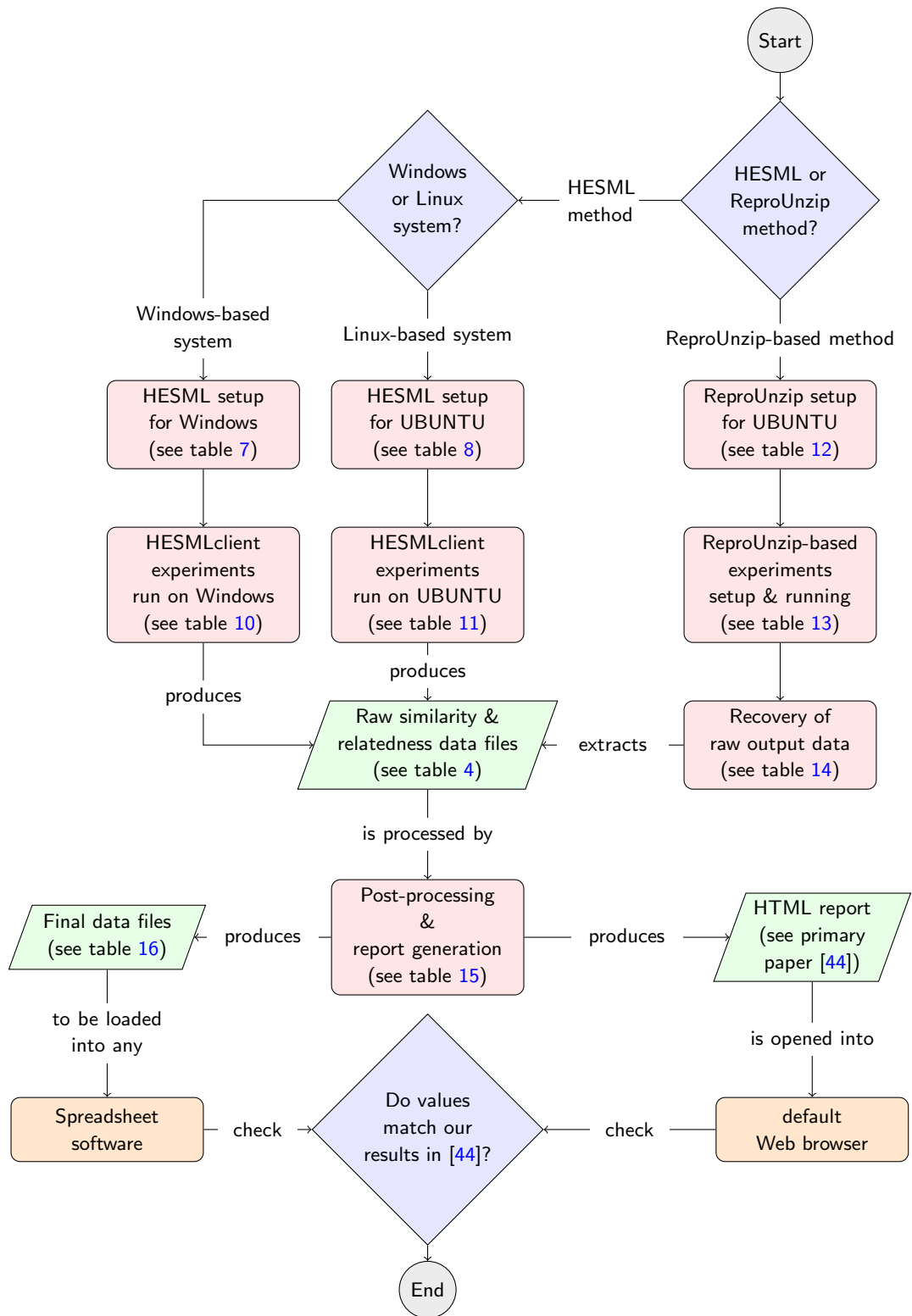


Figure 1: Reproducibility workflows using either HESMLclient or ReprUnzip programs to run the reproducible experiments introduced herein. The three workflows detailed above produce the same raw and processed data files, as well as a collection of HTML pages which reproduce all data tables reported by our primary paper [44].



208 program. Thus, it is enough to download its official distribution from Mende-  
 209 ley [42] or GitHub<sup>9</sup> in order to run our experiments. However, for the sake of  
 210 completeness, section 3.3 introduces the detailed steps to obtain and compile  
 211 HESML V1R4. Finally, we introduce a companion reproducibility dataset [45],  
 212 being publicly available at [46]. This aforementioned reproducibility dataset  
 213 gathers into a common repository all data files required to reproduce our exper-  
 214 iments with the aim of providing a consolidated and permanent version of these  
 215 files, and thus avoiding the tedious work of gathering all these stuff as well as  
 216 any risk of alteration or unavailability of them in the future.

Testing platform	Type	Operating Sys.	Configuration	Tested by
Ubuntu-base1	Virtual	Ubuntu 16.04	1 Core Intel E5-2640-v2 CPU @2 GHz, 8 Gb RAM, 100 Gb SSD disk	Authors
Ubuntu-base2	Virtual	Ubuntu 16.04	2 Intel Core Xeon E5 2699- v4 CPU @2.2 GHz, 8 Gb RAM, 100 Gb SSD disk	Authors
Ubuntu-base3	Virtual	Ubuntu 16.04	1 Core Intel E5-2640-v2 CPU @2 GHz, 8 Gb RAM, 200 Gb SSD disk	Authors
Ubuntu-base4	Virtual	Ubuntu 16.04	2 Intel Core Xeon E5 2699- v4 CPU @2.2 GHz, 8 Gb RAM, 200 Gb SSD disk	Authors
Windows-base1	Laptop	Windows 10x64	1 Intel Core i7-5500U CPU @2.4GHz, 16 Gb RAM, 100 Gb SSD disk	Authors
XXXX-rev1	—	—	—	*Reviewer 1
XXXX-rev2	—	—	—	*Reviewer 2
XXXX-rev3	—	—	—	*Reviewer 3

Table 5: Testing platforms successfully used to reproduce our experiments. Virtual computers are cloud-based servers rented to [clouding.io](https://clouding.io) which are based on the OpenStack virtualization platform [62]. Ubuntu-base1 and Ubuntu-base2, as well as as Ubuntu-base3 and Ubuntu-base4, differ only in the disk space demanded by Reprounzip. On the other hand, Ubuntu-base1 and Ubuntu-base2 differ from Ubuntu-base3 and Ubuntu-base4 in that these two later platforms use a more modern CPU than the former ones, which were used in the implementation of our original experiments in [44]. For this reason, the experiments reproduced on Ubuntu-base3 and Ubuntu-base4 configurations report lower running times than the former ones as shown in table 6. (\*) We will include the testing platforms evaluated by the reviewers.

### 217 3.2. System requirements and performance evaluation

218 Table 5 shows the testing platforms in which we have successfully reproduced  
 219 the experiments detailed herein, whilst table 6 shows their running times in the  
 220 completion of all experiments for each aforementioned reproducibility method  
 221 and testing platform. The configuration of these platforms sets the minimal  
 222 system requirements to reproduce our experiments. Unlike the execution of our  
 223 experiments using *HESMLclient* program on the UBUNTU-based computers

<sup>9</sup><https://github.com/jjlastra/HESML.git>

224 detailed in table 5, the execution using Reprounzip demands more disk space  
 225 because it needs to setup a docker container to run the experiments. For this  
 226 reason, UBUNTU-Reprounzip platforms shown in table 5 are based on a mini-  
 227 mal overall disk space of 200 Gb with the aim of setting up UBUNTU, Docker  
 228 and the resources required by our Reprozip package.

Run	Testing platform	Method	Running time	Tested by
1	Ubuntu-base1	HESMLclient	17581 min $\approx$ 12.2 days	Authors
2	Ubuntu-base3	ReproUnzip	18109 min $\approx$ 12.6 days	Authors
3	Ubuntu-base2	HESMLclient	9622 min $\approx$ 6.68 days	Authors
4	Ubuntu-base4	ReproUnzip	11732 min $\approx$ 8.15 days	Authors
5	Windows-base1	HESMLclient	10 days	Authors
6	Ubuntu-base2	HESMLclient	10201 min $\approx$ 7.08 days	Authors
7				*Reviewer 1
8				*Reviewer 2
9				*Reviewer 3

Table 6: Running times obtained on different testing platforms for the execution of all benchmarks by using HESMLclient program with the 'benchmark\_survey.exp' experiment file, or by running Reprounzip program with the "WN\_ontology\_measures\_vs\_embeddings.rpz" file. (\*) We will include the performance reported by the reviewers.

### 229 3.3. Obtaining and compiling HESML

230 Table 1 shows the technical information required to obtain and compile the  
 231 HESML source code and run the experiments detailed in table 4. HESML V1R4  
 232 distribution includes compiled versions of HESML library and the HESMLclient  
 233 program, thus this later program could be directly used without the need of compiling  
 234 the source code in NetBeans. There are two different ways of obtaining  
 235 the HESML source code as follows: (1) by downloading the latest HESML version  
 236 from the permanent Mendeley Data link [42]; or (2) by downloading it from  
 237 its GitHub repository detailed in table 1. Once the HESML source code has  
 238 been downloaded and extracted onto your hard drive, the project will have the  
 239 folder structure shown in figure 2 and detailed below:

240 **HESML** is the main software library folder containing the NetBeans project  
 241 and HESML source code. Below this folder you find the *dist* folder  
 242 which contains the *HESML-V1R4.jar* distribution file generated during  
 243 the compilation, whilst *HESMLclient* folder contains the source code of the  
 244 HESMLclient console application. The main aim of the *HESMLclient.jar*  
 245 application is to provide a collection of sample functions in order to show  
 246 the HESML functionality, as well as running any (\*.exp) reproducible ex-  
 247 periment file.

248 **PedersenICmodels** folder contains the full WordNet-InfoContent-3.0 collec-  
 249 tion of WordNet-based frequency files created by Ted Pedersen [55]. The  
 250 file names denote the corpus used to build each file. The readme file de-  
 251 tails the method used to build the frequency files, which is also detailed  
 252 in [56].

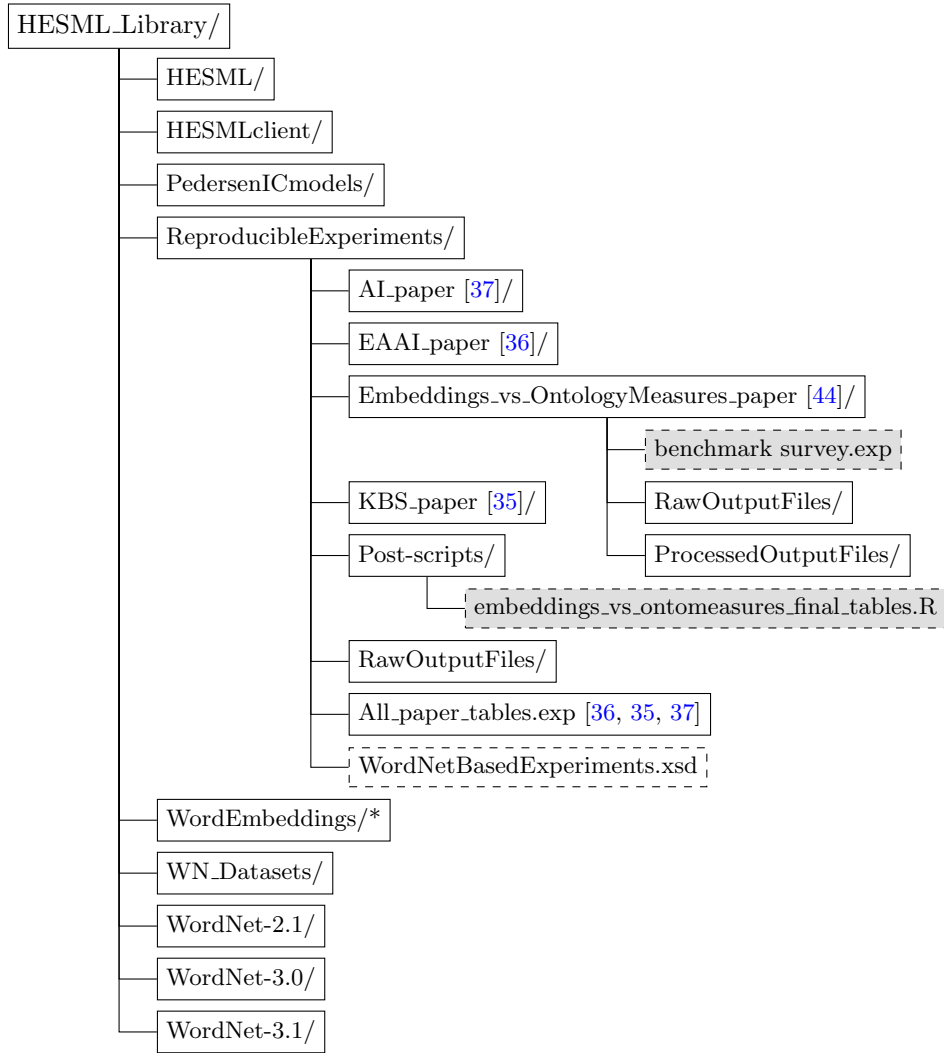


Figure 2: Directory structure of the HESML library once it has been extracted onto disk. The reproducible experiment file and the post-processing R-language script used to reproduce and generate our final data tables respectively are shown in dashed-line boxes in grey, whilst XML-based experiment file format is detailed by XML-schema file shown in unfilled dashed-line box. (\*) WordEmbeddings folder contains the pre-trained files for all word embedding models used in our experiments; however, this folder is neither included by the HESML V1R4 [42] distribution nor HESML V1R4 release at GitHub repository because its large size. Thus, you must download the “*WordEmbeddings.zip*” file [46] and extract it onto the main *HESML\_Library* directory to retrieve this folder and its content.

253 **ReproducibleExperiments** folder contains one subfolder for each paper in-  
254 troduced by Lastra-Díaz and García-Serrano [35, 36, 37] and our pri-  
255 mary paper [44] reproduced herein. Likewise, the aforementioned folder  
256 also contains a XML-schema file called “*WordNetBasedExperiments.xsd*”,  
257 which describes the syntax of all XML-based experiment files (\*.exp),  
258 and the *All\_paper\_tables.exp* file with the definition of all the reproducible  
259 experiments corresponding to the three aforementioned papers of Lastra-  
260 Díaz and García-Serrano. All (\*.exp) files have been created with the  
261 XML Spy editor. In addition, this folder contains the *RawOutputFiles*  
262 subfolder with all the raw output files of the three aforementioned papers  
263 [35, 36, 37].

264 **Post-scripts** folder contains a set of post-processing R scripts which process  
265 the raw output files generated by all reproducible experiments to generate  
266 all final data tables and figures reported in our papers exactly.

267 **WN\_datasets** folder contains a collection of (\*.csv) data files with fields sep-  
268 arated by semicolon which correspond to the word similarity benchmarks  
269 shown in table 4, whilst *WordNet-2.1*, *WordNet-3.0* and *WordNet-3.1*  
270 contain the database files of three different versions of WordNet.

271 **Embeddings\_vs\_OntologyMeasures\_paper** folder contains the reproducible  
272 experiment file “benchmark\_survey.exp” encoding all benchmarks intro-  
273 duced herein and detailed in table 4. In addition, this folder contains  
274 a subfolder called “RawOutputFiles” containing all raw output similar-  
275 ity files generated by our experiments. The R-language script file called  
276 “embeddings\_vs\_ontomeasures\_final\_tables.R” generates all files in “Pro-  
277 cessedOutputFiles” subfolder.

278 Tables 7 and 8 show a detailed step-by-step procedure to set up our repro-  
279 ducible experiments based on HESML on any Windows or Linux-based system  
280 respectively. HESML distribution includes pre-compiled versions of HESML-  
281 client program and HESML library; thus, you could skip the compilation step  
282 for running our experiments. However, for the sake of completeness, we briefly  
283 detail the compilation steps in table 9.

Step	Windows-based setup instructions for HESMLclient experiments
(1)	Install Java 8 runtime or higher in your workstation.
(2)	Open a PowerShell console (Windows 7 and higher) in any directory.
(3)	Create a working directory and move to it as follows: <pre>\$ mkdir REPRODIR \$ cd REPRODIR</pre>
(4)	Download and extract latest HESML version from its GitHub repository (see URL below) using either any Web browser or PowerShell as detailed below: <pre>\$ powershell -command "&amp; { iwr https://github.com/jjlastra/HESML/archive/master.zip }" \$ Expand-Archive ./master.zip .</pre>
(5)	Download the <i>WordEmbeddings.zip</i> file from our Dataverse repository [46, see URL below] and extract it onto HESML root directory using either any Web browser or PowerShell as detailed below: <pre>\$ cd HESMLmaster/HESML_LIBRARY \$ mkdir WordEmbeddings \$ cd WordEmbeddings \$ powershell -command "&amp; { iwr https://doi.org/10.21950/wordembeddings.zip }" \$ Expand-Archive ./wordembeddings.zip .</pre>

Table 7: Detailed instructions for downloading HESML V1R4 onto a Windows-based system from its GitHub repository.

Step	Linux-based setup instructions for HESMLclient experiments
(1)	Install Java 8 and Java SE Dev Kit 8 or higher as follows: <pre>\$ sudo apt-get update \$ sudo apt-get -y install default-jdk</pre>
(2)	Install UNZIP program as follows: <pre>\$ sudo apt-get update \$ sudo apt-get -y install unzip</pre>
(3)	Create a working directory and move to it as follows: <pre>\$ mkdir REPRODIR \$ cd REPRODIR</pre>
(4)	Download the latest HESML version from GitHub (see URL below) as follows: <pre>\$ wget https://github.com/jjlastra/HESML/archive/master.zip</pre>
(5)	Extract <i>master.zip</i> file onto your working directory as follows: <pre>\$ unzip master.zip</pre>
(6)	Download the <i>WordEmbeddings.zip</i> file from our Dataverse repository [46, see URL below] and extract it onto HESML root directory as detailed below. <pre>\$ cd HESMLmaster/HESML_LIBRARY \$ wget https://doi.org/10.21950/wordembeddings.zip \$ unzip wordembeddings.zip -d WordEmbeddings</pre>

Table 8: Detailed instructions for downloading HESML V1R4 from its GitHub repository onto a Linux-based system.

Step	Detailed instructions to compile HESML
------	--

- |     |   |
|-----|---|
| (1) | Follow the step-by-step procedures to download HESML source code as detailed in tables 7 and 8 for Windows or Linux-based systems respectively.   |
| (2) | Install Java SE Dev Kit 8 and NetBeans 8.0.2 or higher in your workstation.   |
| (3) | Launch NetBeans IDE and open the <i>HESML</i> and <i>HESMLclient</i> projects contained in the HESML root folder as shown in table 2. NetBeans automatically detects the presence of a <i>nbproject</i> subfolder with the project files. |
| (4) | Select <i>HESML</i> and <i>HESMLclient</i> projects in the project treeview respectively. Then, invoke the “Clean and Build project (Shift + F11)” command in order to compile both projects.   |

Table 9: Detailed instructions for compiling HESML onto any Windows or Linux-based system. We recall that the compilation of HESML is unneeded to run all experiments introduced herein.

### 284 3.4. Running the experiments with HESMLclient

285 Once you have downloaded and extracted the *HESML VIR4* library onto  
 286 your hard drive as detailed in section 3.3, you are ready to run the reproducible  
 287 experiments by following the steps detailed in tables 10 and 11 for testing plat-  
 288 forms based on Windows and Linux respectively. However, before to run the  
 289 experiments, you must download the *WordEmbeddings.zip* file [46] and extract  
 290 it onto the main *HESML\_Library* directory as detailed in step 5 of table 7 for  
 291 Windows, and step 6 of table 8 for the Linux-based case. This later ZIP file  
 292 contains all pre-trained word embedding files; however, it is not included in  
 293 current HESML distribution because of its large size and the space limitations  
 294 of GitHub and Mendeley repositories. We note that the original *HESMLclient*  
 295 source code is defined to fetch the required input files from the folder structure  
 296 of *HESML* as shown in figure 2.

Step	HESMLclient running instructions on any Windows-based system
------	--

- |     |   |
|-----|---|
| (1) | Open a command console in the <i>HESMLclient</i> directory as shown in figure 2.<br><code>\$ cd REPROD/DIR/HESML_Library/HESMLclient</code>   |
| (2) | Run the following command with the reproducible experiment file:<br><br><code>\$ java -jar -Xms4096m dist/HESMLclient.jar ../ReproducibleExperiments/Embeddings_vs_OntologyMeasures_paper/benchmark_survey.exp</code> |
| (3) | Command in step 2 above will generate all raw output files listed in table 4 onto <code>../ReproducibleExperiments/Embeddings_vs_OntologyMeasures_paper</code> folder (see figure 2).                                 |

Table 10: Detailed instructions for running our experiments with the *HESMLclient* program on any testing platform based on Windows.

Step	HESMLclient running instructions on any Linux-based system
(1)	Open a Linux command console in the <i>HESMLclient</i> directory (see figure 2). <pre>user@server\$ cd REPRODIR/HESMLLibrary/HESMLclient</pre>
(2)	We create a “screen” session and run HESMLclient in background. Note that HESMLclient execution could take up to two weeks (see table 6). <pre>user@server\$ screen -S REPROEXPS user@screen\$ java -jar -Xms4096m dist/HESMLclient.jar ../ReproducibleExperiments/Embeddings-vs-OntologyMeasures-paper/benchmark_survey.exp</pre>
(3)	We detach from “screen” before to close the server main console <pre>user@screen\$ CTRL+a, d</pre>
(4)	We reattach to the screen console to check the completion of HESMLclient <pre>user@server\$ screen -r REPROEXPS</pre>
(5)	We destroy the “screen” console once finished HESMLclient execution <pre>user@server\$ screen -X -S REPROEXPS quit</pre>
(6)	Second command in step (2) above will generate all raw output files listed in table 4 onto <i>../ReproducibleExperiments/Embeddings-vs-OntologyMeasures-paper</i> folder (see figure 2).

Table 11: Detailed instructions for running our experiments with the *HESMLclient* program on any testing platform based on Linux.

### 297 3.5. Running the *ReproZip* experiments

298 The *ReproZip*<sup>10</sup> program was used for recording and packaging the running  
299 of the *HESMLclient* program with all the reproducible experiments defined by  
300 the “*benchmark\_survey.exp*” file into the “*WN\_ontology\_measures\_vs\_embeddings.rpz*”  
301 file, which is publicly available at our UNED Dataverse repository [46]. This  
302 later *Reprozip* file was generated by running *Reprozip* on the Ubuntu-base1  
303 workstation detailed in table 5; however, in order to run *ReproUnzip* based on  
304 Docker as detailed below is needed to set up an Ubuntu-*Reprounzip* platform  
305 (see table 5). Because the execution of the experiments takes long time, and  
306 *Reprounzip* with Docker cannot be executed in background mode without any  
307 output console, we will setup and use “screen” program on Linux.

308 In order to set up and run the reproducible experiments introduced herein,  
309 you need to use *ReproUnzip*. *ReproUnzip* can be used with two different virtu-  
310 alization platforms: (1) Vagrant + VirtualBox, or (2) Docker. However, because  
311 of its simple setup and computational efficiency, our preferred *ReproUnzip* con-  
312 figuration is that based on Docker. For instance, in order to setup *ReproUnzip*  
313 based on Docker for Ubuntu, you should follow the detailed steps shown in table  
314 12, despite several steps possibly being unnecessary depending on your starting  
315 configuration. Once *ReproUnzip* and Docker have been successfully installed,  
316 table 13 shows the detailed instructions to set up and run the reproducible  
317 experiments. Those readers who prefer to use *ReproUnzip* with VirtualBox  
318 instead of Docker can consult the *ReproZip* installation page<sup>11</sup>.

319 The running of the reproducible experiments based on Docker for Ubuntu  
320 tooks approximately one week in a modern virtual computer as detailed in  
321 table 6. Once the running is completed, you should follow the instructions

<sup>10</sup><https://www.reprozip.org/>

<sup>11</sup><https://reprozip.readthedocs.io/en/1.0.x/install.html>

Steps (1-4) below install Reprounzip and all its dependencies.

- (1) `$ sudo apt-get update`
- (2) `$ sudo apt-get -y install libffi-dev libssl-dev openssl openssh-server`
- (3) `$ sudo apt-get -y install libsqlite3-dev python-dev python-pip screen`
- (4) `$ sudo pip install reprounzip[all]`

Steps (5-11) below install latest version of Docker CE whilst step 11 checks its installation. For further details, we refer the reader to the official Docker setup page: <https://docs.docker.com/install/linux/docker-ce/ubuntu/>

- (5) `$ sudo apt-get -y install apt-transport-https ca-certificates`
- (6) `$ sudo apt-get -y install curl gnupg-agent software-properties-common`
- (7) `$ curl -fsSL https://download.docker.com/linux/ubuntu/gpg | sudo apt-key add -`
- (8) `$ sudo add-apt-repository "deb [arch=amd64] https://download.docker.com/linux/ubuntu $(lsb_release -cs) stable"`
- (9) `$ sudo apt-get update`
- (10) `$ sudo apt-get -y install docker-ce docker-ce-cli containerd.io`
- (11) `$ sudo docker run hello-world`

Table 12: Detailed instructions on installing ReproUnzip with Docker for Ubuntu. Despite that steps above could look tedious, we prefer that readers are aware of all packages being installed instead of running a single setup script hiding this information.

322 shown in table 14 to retrieve the raw output files from the Docker container,  
323 as listed in table 4. Finally, table 6 reports a sample of software platforms in  
324 which the Reprozip-based experiments introduced herein have been successfully  
325 reproduced.

### 326 3.6. Processing of the raw output files

327 The running of the “*benchmark\_survey.exp*” experiment file generates the  
328 collection of comma-separated files (\*.csv) listed in table 4, whose values are  
329 separated by a semicolon. All raw output files are saved in the same folder as  
330 their corresponding input reproducible experiment file.

331 Raw output similarity files generated by our experiments must be processed  
332 in order to compute the Pearson, Spearman and Harmonic score metrics match-  
333 ing the tables shown in our primary paper [44]. We provide a R-language script  
334 called “*embeddings\_vs\_ontomeasures\_final\_tables.R*” with the aim of automating  
335 this post-processing. Latest version of the aforementioned post-processing script  
336 should be obtained from HESML GitHub distribution, as detailed in tables 7  
337 and 8, or from our companion reproducibility dataset [46]. This aforementioned  
338 script includes the source code of the *mat.sort* function provided by the Bio-  
339 PhysConnectoR package [31], which is no longer available in CRAN server.

340 In order to carry-out the aforementioned post-processing, you should setup  
341 the well-known R statistical program<sup>12</sup> in your workstation and follow the steps

---

<sup>12</sup><https://www.r-project.org/>



- (1) Setup the Rezip program onto any supported platform (Linux, Windows and MacOS) by following the step-by-step guide detailed in table 12 (see Rezip installation page for further information).
- (2) Move to the home directory and create a working directory as follows
 

```
$ cd /home
$ mkdir REPROEXPS
$ cd REPROEXPS
```
- (3) Download the *WN\_ontology\_measures\_vs\_embeddings.rpz* (12.4 Gb) from its repository [46]. For instance, you can execute the command below. The download of this file could takes several minutes.
 

```
$ wget https://doi.org/10.21950/wn_ontology_measures_vs_embeddings.rpz
```
- (4) Next, we must setup the docker container as detailed below which could take up to 45 minutes depending of your testing platform. Thus, we recommend to create a “screen” session to run in background both setup and running of the Rezip-based experiment. You can detach from “screen” console by pressing “Ctrl+a,d”.
 

```
user@server$ screen -S REPROUNZIP
user@server$ rezip docker setup
wn_ontology_measures_vs_embeddings.rpz docker_folder
```
- (5) Next, we will run the Rezip-based experiment. Note that Rezip execution could take up to two weeks depending on your hardware setup (see table 6). We strongly recommend to keep open the *screen* console to run the experiment in background as detailed below.
 

```
user@screen$ rezip docker run docker_folder
```

  - We detach from “screen” before to close the server main console
 

```
user@screen$ CTRL+a, d
```
  - We reattach to the screen console to check the completion of Rezip
 

```
user@server$ screen -r REPROUNZIP
```
  - We destroy the “screen” console once finished Rezip execution
 

```
user@server$ screen -X -S REPROUNZIP quit
```

Table 13: Detailed instructions on how to reproduce the packaged experiments once Rezip has been installed. We use *screen* program with the aim of allowing the execution of Rezip in background whilst main program console is detached and closed.

342 detailed in table 15. Then, you need to install the “knitr” and “readr” packages  
 343 using the functionality provided for this task by R program. Table 16 shows the  
 344 output files which are generated from the raw output files listed in table 4 by  
 345 running our aforementioned post-processing script, as well as their correspond-  
 346 ing data tables in our primary paper [44]. In addition, our post-processing script  
 347 generates a collection of HTML files which contain all data tables reported in  
 348 our primary paper [44].

349 Finally, raw data files and processed data files shown in tables 4 and 16  
 350 respectively could be loaded into any spreadsheet software to carry-out any  
 351 further data analysis or confirming the reproducibility of the experiments and  
 352 results reported by our primary paper [44].

---

Step	Recovering the output files generated by our Reprozip-based experiments
------	---

---

- |   |   |
|---|---|
| 1 | user@server\$ reprozip showfiles docker_folder                  |
| 2 | user@server\$ sudo reprozip docker download --all docker_folder |

Table 14: Detailed instruction to recover the output files generated by our Reprozip-based experiments. The first instruction shows a list with the output files generated by the experiments, whilst the second one extracts all the output files from the container and downloads them onto the current folder. You should obtain all raw output files listed in table 4.

---

Step	Detailed post-processing instructions based on a R-language script
------	--

---

- |     |  |
|-----|--|
| (1) | Launch the R statistical program and install <i>knitr</i> and <i>readr</i> packages.   |
| (2) | Launch the R statistical program (you could also use R-Studio).  |
| (3) | Select the menu option “ <i>File-&gt;Open script</i> ”. Then, load the R-language script file called <i>embeddings_vs_ontomeasures_final_tables.R</i> contained in the folder shown in figure 2. Latest version of the aforementioned script should be obtained from HESML GitHub distribution or our companion reproducibility dataset [46].  |
| (4) | Edit the <i>rootDir</i> , <i>inputDir</i> and <i>outputDir</i> variables at the beginning of the script in order to set the directory which contains the raw output files onto your hard drive, as well as the directory in which will be saved the final assembled data tables as reported in our primary paper [44]. IMPORTANT NOTE: <i>inputDir</i> and <i>outputDir</i> variables should end with slash <i>’/’</i> symbol. |
| (5) | Select the menu option “ <i>Edit-&gt;Run all</i> ”. The final assembled data tables will be saved in the output directories defined above, as detailed in table 16. In addition, the aforementioned R script creates a and opens a collection of HTML files which show all data tables in our primary paper [44] and detailed in table 16.   |

Table 15: Detailed instructions for the post-processing of the raw output files generated by our experiments. R-language script computes all Pearson, Spearman and Harmonic score metrics and generates a HTML report file reproducing all data tables reported by our primary paper Lastra-Díaz et al. [44].

#### 353 4. Extending and reusing our reproducible experiments

354 Our reproducible experiments are based on the XML-based HESML ex-  
355 periment file format (\*.exp) whose specification is detailed by the “*WordNet-*  
356 *BasedExperiments.xsd*” schema file distributed with HESML library as shown  
357 in figure 2. Both \*.exp experiment files and \*.xsd schema file were created with  
358 XML Spy editor. Next paragraphs provide a detailed description of the main  
359 objects encoded by the HESML XML-based experiment file format and how  
360 they could be used to create new experiments from scratch as those introduced  
361 herein.

362 *HESML XML-based experiment file format.* Figure 3 shows a sample file which  
363 has been extracted from the “benchmark\_survey.exp” file encoding all repro-  
364 ducible experiments introduced herein. *WordNetBasedExperiments* is the root  
365 node which contains the collection of word similarity or relatedness benchmarks  
366 to be evaluated, whilst *SingleDatasetSimilarityValuesExperiment* encodes a spe-  
367 cific word similarity benchmark defined by a dataset, an output directory and a  
368 collection of WordNet-based similarity measures and pre-trained word embed-

File#	Post-processing output files saved at “outputDir” directory	In primary paper [44]
1	table_Pearson_SimDatasets.csv	table 4 (full precision)
2	table_Pearson_SimDatasets_rounded.csv	table 4
3	table_Spearman_SimDatasets.csv	table 5 (full precision)
4	table_Spearman_SimDatasets_rounded.csv	table 5
5	table_Pearson_RelDatasets.csv	table 6 (full precision)
6	table_Pearson_RelDatasets_rounded.csv	table 6
7	table_Spearman_RelDatasets.csv	table 7 (full precision)
8	table_Spearman_RelDatasets_rounded.csv	table 7
9	table_joined_allEmbeddings_similarity.csv	table 8 (full precision)
10	table_joined_allEmbeddings_similarity_rounded.csv	table 8
11	table_joined_allEmbeddings_relatedness.csv	table 9 (full precision)
12	table_joined_allEmbeddings_relatedness_rounded.csv	table 9
13	table_pvalues_AttractReppel_allembeddings_similarity.csv	table A.1
14	table_pvalues_Paragramws_allembeddings_relatedness.csv	table A.2
15	table_AvgMeasures_Pearson_SimDatasets.csv	table A.3 (full precision)
16	table_AvgMeasures_Pearson_SimDatasets_rounded.csv	table A.3
17	table_AvgMeasures_Spearman_SimDatasets.csv	table A.4 (full precision)
18	table_AvgMeasures_Spearman_SimDatasets_rounded.csv	table A.4
19	table_AvgMeasures_Pearson_RelDatasets.csv	table A.5 (full precision)
20	table_AvgMeasures_Pearson_RelDatasets_rounded.csv	table A.5
21	table_AvgMeasures_Spearman_RelDatasets.csv	table A.6 (full precision)
22	table_AvgMeasures_Spearman_RelDatasets_rounded.csv	table A.6

Table 16: Collection of processed output files generated by the execution of the “embeddings\_vs\_ontomeasures\_final\_tables.R” script file onto the *outputDir* directory and their corresponding tables in our primary work [44].

369 ding models. *SimilarityMeasures* nodes encode ontology-based semantic simi-  
370 larity measures based on WordNet which could require a further Information  
371 Content (IC) model for its implementation, being declared below the *Word-*  
372 *NetMeasures* node. Likewise, *RawWordVectorFiles* encodes the collection of  
373 pre-trained word embedding files to be evaluated in the same dataset. Both  
374 *SimilarityMeasures* and *RawWordVectorFiles* could be declared independently,  
375 and they could contain an unlimited number of methods to be evaluated. Latest  
376 HESML version supports three different pre-trained word embeddings file for-  
377 mats which are defined by the *EmbVectorFiles*, *UKBVectorFiles* and *NasariVec-*  
378 *torFiles* nodes. Raw output files which are generated by *SingleDatasetSimilar-*  
379 *ityValuesExperiment* procedures contain a matrix of values encoding the raw  
380 similarity value reported by each method for each word pair in the similarity  
381 dataset being evaluated.

## XML-based HESML experiment file sample

---

```
<WordNetBasedExperiments>
  <SingleDatasetSimilarityValuesExperiment>
    <OutputFileName>raw_similarity_values_MC28_dataset.csv</OutputFileName>
    <DatasetDirectory>../WN\_Datasets</DatasetDirectory>
    <DatasetFileName>Miller_Charles_28_dataset.csv</DatasetFileName>
    <WordNetMeasures>
      <WordNetDatabaseFileName>data.noun</WordNetDatabaseFileName>
      <WordNetDatabaseDirectory>../Wordnet-3.0/dict</WordNetDatabaseDirectory>
      <SimilarityMeasures>
        <SpecificSimilarityMeasure>
          <MeasureType>JiangConrath</MeasureType>
          <IntrinsicICModel>Sanchez2011</IntrinsicICModel>
        </SpecificSimilarityMeasure>
        <SpecificSimilarityMeasure>
          <MeasureType>Rada</MeasureType>
        </SpecificSimilarityMeasure>
      </SimilarityMeasures>
    </WordNetMeasures>
    <RawWordVectorFiles>
      <EmbVectorFiles>
        <VectorFile>../WordEmbeddings/attract-reppel.emb</VectorFile>
      </EmbVectorFiles>
      <UKBVectorFiles>
        <VectorFile>../WordEmbeddings/wordnet-ukb.ppv</VectorFile>
      </UKBVectorFiles>
      <NasariVectorFiles>
        <NasariVectorFile>
          <SensesFile>../WordEmbeddings/nasari/en_wordsenses_BN.txt</SensesFile>
          <VectorFile>../WordEmbeddings/nasari/nasari-unified</VectorFile>
        </NasariVectorFile>
      </NasariVectorFiles>
    </RawWordVectorFiles>
  </SingleDatasetSimilarityValuesExperiment>
</WordNetBasedExperiments>
```

Figure 3: XML source code above shows an example of a HESML reproducible experiment on word similarity and relatedness. Source code above has been extracted from the “benchmark\_survey.exp” file which encodes all experiments reported in our primary paper [44].

382 *Extending or modifying our experiments.* Anyone could use our main aforemen-  
383 tioned experiment file as a template to create new experiments from scratch  
384 by evaluating other sets of available ontology-based semantic similarity mea-  
385 sures based on WordNet, pre-trained word embedding models or word similar-  
386 ity datasets not considered herein. For instance, any reader could evaluate any  
387 ontology-based methods by declaring it in any *SimilarityMeasures* node when-  
388 ever this method have been previously implemented in HESML and its keyname  
389 being specified by the *SimilarityMeasureType* enumeration in the “*WordNet-*  
390 *BasedExperiments.xsd*” schema file. Likewise, currently supported IC models  
391 are specified by the *IntrinsicICModelType* enumeration in the aforementioned  
392 XML schema file. On the other hand, any reader could evaluate any uncon-  
393 sidered pre-trained word embedding model by declaring a new method in the  
394 *RawWordVectorFiles*, whenever its corresponding pre-trained model being pro-  
395 vided in any of the three file formats which are currently supported, otherwise  
396 there would be needed to extend HESML to support a new pre-trained word  
397 embedding file format. Finally, any reader could define any new benchmark  
398 considering a different set of word similarity datasets by declaring further *Sin-*  
399 *gleDatasetSimilarityValuesExperiment* nodes with their corresponding dataset  
400 files in comma-separated file format. For a detailed list of the methods currently  
401 implemented by HESML V1R4, we refer the readers to its release notes [42].

## 402 5. Conclusions and future work

403 This work introduces for the first time a large set of reproducible experi-  
404 ments on word similarity and relatedness including most methods in the fami-  
405 lies of ontology-based semantic similarity measures based on WordNet and word  
406 embedding models. This aforementioned set of experiments allow that all ex-  
407 periments and results introduced by Lastra-Díaz et al. [44] to be reproduced  
408 exactly. Likewise, our reproducible experiments could be easily extended or  
409 modified to create new benchmarks from scratch which evaluate a different set  
410 of methods and word similarity and relatedness datasets from those considered  
411 herein. For this reason, we hope that this set of reproducible benchmarks to be-  
412 come into a de facto standard experimentation platform for any future research  
413 on word similarity and relatedness.

414 As forthcoming activities, we plan the study and proposal of new distribu-  
415 tional similarity and relatedness measures, as well as their use in the definition  
416 of sentence and short-text similarity measures.

## 417 Acknowledgements

418 We are grateful of Fernando González and Juan Corrales for setting up our  
419 UNED Dataverse dataset, Yuanyuan Cai for answering kindly our questions  
420 to replicate their IC-based similarity measures and IC models in HESML, and  
421 <http://clouding.io> for their technical support to set up our experimental  
422 platform. We are also very thankful to José Camacho-Collados for providing  
423 the weighting overlap source code which we have integrated into HESML for  
424 measuring the similarity between the NASARI vectors. This work has been  
425 partially supported by the Spanish project VEMODALEN (TIN2015-71785-R),

426 the Basque Government (type A IT1343-19), BBVA BigKnowledge bigknowl-  
427 edge project and the Spanish Research Agency LIHLITH project (PCIN-2017-  
428 118 / AEI) in the framework of EU ERA-Net CHIST-ERA.

## 429 References

- 430 [1] Agirre, E., Alfonseca, E., Hall, K., Kravalova, J., Paşca, M., Soroa, A.,  
431 2009. A Study on Similarity and Relatedness Using Distributional and  
432 WordNet-based Approaches. In: Proc. of Human Language Technolo-  
433 gies: The 2009 Annual Conf. of the North American Chapter of the  
434 Association for Computational Linguistics. NAACL '09. ACL, Strouds-  
435 burg, PA, USA, pp. 19–27.
- 436 [2] Agirre, E., Soroa, A., 2009. Personalizing pagerank for word sense disam-  
437 biguation. In: Proc. of the 12th Conf. of the European Chapter of the  
438 Association for Computational Linguistics. ACL, pp. 33–41.
- 439 [3] Ben Aouicha, M., Hadj Taieb, M. A., Nov. 2015. G2WS: Gloss-based Word-  
440 Net and Wiktionary semantic Similarity measure. In: 2015 IEEE/ACS  
441 12th International Conference of Computer Systems and Applications  
442 (AICCSA). pp. 1–7.
- 443 [4] Ben Aouicha, M., Hadj Taieb, M. A., 2016. Computing semantic similarity  
444 between biomedical concepts using new information content approach.  
445 Journal of Biomedical Informatics 59, 258–275.
- 446 [5] Ben Aouicha, M., Hadj Taieb, M. A., Ben Hamadou, A., Nov. 2016. SISR:  
447 System for integrating semantic relatedness and similarity measures.  
448 Soft Computing, 1–25.
- 449 [6] Ben Aouicha, M., Hadj Taieb, M. A., Ben Hamadou, A., Mar.  
450 2016. Taxonomy-based information content and wordnet-wiktionary-  
451 wikipedia glosses for semantic relatedness. Applied Intelligence, 1–37.
- 452 [7] Ben Aouicha, M., Hadj Taieb, M. A., Beyaoui, S., Aug. 2016. Distributional  
453 semantics study using the co-occurrence computed from collaborative  
454 resources and WordNet. In: 2016 International Symposium on INnova-  
455 tions in Intelligent SysTems and Applications (INISTA). pp. 1–8.
- 456 [8] Bojanowski, P., Grave, E., Joulin, A., Mikolov, T., Jul. 2016. Enriching  
457 Word Vectors with Subword Information. arXiv:1607.04606.
- 458 [9] Branco, A., Cohen, K. B., Vossen, P., Ide, N., Calzolari, N., Mar. 2017.  
459 Replicability and reproducibility of research results for human language  
460 technology: introducing an LRE special section. Language Resources  
461 and Evaluation 51 (1), 1–5.
- 462 [10] Bruni, E., Tran, N.-K., Baroni, M., 2014. Multimodal Distributional Se-  
463 mantics. Journal of Artificial Intelligence Research 49 (1), 1–47.
- 464 [11] Camacho-Collados, J., Pilehvar, M. T., Navigli, R., 2016. Nasari: Integrat-  
465 ing explicit knowledge and corpus statistics for a multilingual represen-  
466 tation of concepts and entities. Artif. Intell. 240, 36–64.

- 467 [12] Chirigati, F., Capone, R., Rampin, R., Freire, J., Shasha, D., Mar. 2016.  
468 A collaborative approach to computational reproducibility. *Information*  
469 *Systems* 59, 95–97.
- 470 [13] Chirigati, F., Rampin, R., Shasha, D., Freire, J., 2016. ReProZip: compu-  
471 tational reproducibility with ease. In: *Proc. of the ACM Intl. Conf. on*  
472 *Management of Data (SIGMOD)*. Vol. 16. pp. 2085–2088.
- 473 [14] Fariña, A., Martínez-Prieto, M. A., Claude, F., Navarro, G., Lastra-Díaz,  
474 J. J., Prezza, N., Seco, D., Apr. 2019. On the reproducibility of experi-  
475 ments of indexing repetitive document collections. *Information systems*  
476 83, 181–194.
- 477 [15] Finkelstein, L., Gabrilovich, E., Matias, Y., Rivlin, E., Solan, Z., Wolfman,  
478 G., Ruppin., E., 2002. Placing search in context: the concept revisited.  
479 *ACM Trans. on Information Systems* 20 (1), 116–131.
- 480 [16] Fokkens, A., Van Erp, M., Postma, M., Pedersen, T., Vossen, P., Freire, N.,  
481 Aug. 2013. Offspring from Reproduction Problems: What Replication  
482 Failure Teaches Us. In: *Proceedings of the 51st Annual Meeting of the*  
483 *Association for Computational Linguistics. ACL, Sofia, Bulgaria*, pp.  
484 1691–1701.
- 485 [17] Gerz, D., Vulić, I., Hill, F., Reichart, R., Korhonen, A., Nov. 2016.  
486 SimVerb-3500: A Large-Scale Evaluation Set of Verb Similarity. In:  
487 *Proc. of EMNLP. Austin, Texas*, pp. 2173–2182.
- 488 [18] Goikoetxea, J., Agirre, E., Soroa, A., 2014. Exploring the Use of Word  
489 Embeddings and Random Walks on Wikipedia for the CogAlex Shared  
490 Task. *Proc. of the 4th Workshop on Cognitive Aspects of the Lexicon*,  
491 31–34.
- 492 [19] Goikoetxea, J., Agirre, E., Soroa, A., 2016. Single or Multiple? Combining  
493 Word Representations Independently Learned from Text and WordNet.  
494 In: *Proc. of AAAI*. pp. 2608–2614.
- 495 [20] Goikoetxea, J., Soroa, A., Agirre, E., 2015. Random Walks and Neural  
496 Network Language Models on Knowledge Bases. In: *Proc. of HLT-*  
497 *NAACL*. pp. 1434–1439.
- 498 [21] Goikoetxea, J., Soroa, A., Agirre, E., Jul. 2018. Bilingual Embeddings with  
499 Random Walks over Multilingual Wordnets. *Knowledge-Based Systems*  
500 150 (15), 218–230.
- 501 [22] Hadj Taieb, M. A., Ben Aouicha, M., Ben Hamadou, A., Nov. 2014. A new  
502 semantic relatedness measurement using WordNet features. *Knowledge*  
503 *and Information Systems* 41 (2), 467–497.
- 504 [23] Hadj Taieb, M. A., Ben Aouicha, M., Ben Hamadou, A., Nov. 2014.  
505 Ontology-based approach for measuring semantic similarity. *Engineering*  
506 *Applications of Artificial Intelligence* 36, 238–261.

- 507 [24] Hadj Taieb, M. A., Ben Aouicha, M., Bourouis, Y., Jun. 2015. FM3S:  
508 Features-Based Measure of Sentences Semantic Similarity. In: Onieva,  
509 E., Santos, I., Osaba, E., Quintián, H., Corchado, E. (Eds.), Proceed-  
510 ings of the 10th International Conference on Hybrid Artificial Intelligent  
511 Systems (HAIS 2015). Vol. 9121 of LNCS. Springer, Bilbao, Spain, pp.  
512 515–529.
- 513 [25] Hadj Taieb, M. A., Ben Aouicha, M., Tmar, M., Ben Hamadou, A., 2012.  
514 New WordNet-based semantic relatedness measurement. In: GWC 2012  
515 6th International Global Wordnet Conference. academia.edu, p. 126.
- 516 [26] Hadj Taieb, M. A., Ben Aouicha, M., Tmar, M., Ben Hamadou, A., Nov.  
517 2012. Wikipedia Category Graph and New Intrinsic Information Con-  
518 tent Metric for Word Semantic Relatedness Measuring. In: Proc. of the  
519 Third International Conference on Data and Knowledge Engineering  
520 (ICDKE). Lecture Notes in Computer Science. Springer Berlin Heidel-  
521 berg, Wuyishan, Fujian, China, pp. 128–140.
- 522 [27] Hadj Taieb, M. A., Ben Aouicha, M., Tmar, M., Hamadou, A. B., 2011.  
523 New information content metric and nominalization relation for a new  
524 WordNet-based method to measure the semantic relatedness. In: Cy-  
525 bernetetic Intelligent Systems (CIS), 2011 IEEE 10th International Con-  
526 ference on. pp. 51–58.
- 527 [28] Halawi, G., Dror, G., Gabrilovich, E., Koren, Y., 2012. Large-scale Learn-  
528 ing of Word Relatedness with Constraints. In: Proc. of ACM SIGKDD.  
529 ACM, New York, NY, USA, pp. 1406–1414.
- 530 [29] Harispe, S., Ranwez, S., Janaqi, S., Montmain, J., Mar. 2014. The semantic  
531 measures library and toolkit: fast computation of semantic similarity  
532 and relatedness using biomedical ontologies. *Bioinformatics* 30 (5), 740–  
533 742.
- 534 [30] Hill, F., Reichart, R., Korhonen, A., Dec. 2015. SimLex-999: Evaluating  
535 Semantic Models with (Genuine) Similarity Estimation. *Computational*  
536 *Linguistics* 41 (4), 665–695.
- 537 [31] Hoffgaard, F., Weil, P., Hamacher, K., Apr. 2010. BioPhysConnectoR: Con-  
538 necting sequence information and biophysical models. *BMC bioinforma-*  
539 *tics* 11, 199.
- 540 [32] Huang, E. H., Socher, R., Manning, C. D., Ng, A. Y., 2012. Improving Word  
541 Representations via Global Context and Multiple Word Prototypes. In:  
542 Proc. of the Annual Meeting of the ACL. Vol. 1. pp. 873–882.
- 543 [33] Lastra-Díaz, J. J., Sep. 2014. Intrinsic Semantic Spaces for the representa-  
544 tion of documents and semantic annotated data. Master’s thesis, Uni-  
545 versidad Nacional de Educación a Distancia (UNED). Department of  
546 Computer Languages and Systems, [http://e-spacio.uned.es/fez/  
547 view/bibliuned:master-ETSIInformatica-LSI-Jlastra](http://e-spacio.uned.es/fez/view/bibliuned:master-ETSIInformatica-LSI-Jlastra).
- 548 [34] Lastra-Díaz, J. J., Sep. 2017. Recent Advances in Ontology-based Se-  
549 mantic Similarity Measures and Information Content Models based



- 550 on WordNet. Ph.D. thesis, Universidad Nacional de Educación a Dis-  
551 tancia (UNED), [http://e-spacio.uned.es/fez/view/tesisuned:](http://e-spacio.uned.es/fez/view/tesisuned:ED-Pg-SisInt-Jjlastra)  
552 [ED-Pg-SisInt-Jjlastra](http://e-spacio.uned.es/fez/view/tesisuned:ED-Pg-SisInt-Jjlastra).
- 553 [35] Lastra-Díaz, J. J., García-Serrano, A., Nov. 2015. A new family of in-  
554 formation content models with an experimental survey on WordNet.  
555 Knowledge-Based Systems 89, 509–526.
- 556 [36] Lastra-Díaz, J. J., García-Serrano, A., Nov. 2015. A novel family of IC-  
557 based similarity measures with a detailed experimental survey on Word-  
558 Net. Eng. App. of Artif. Intell. 46, 140–153.
- 559 [37] Lastra-Díaz, J. J., García-Serrano, A., Jul. 2016. A refinement of  
560 the well-founded Information Content models with a very de-  
561 tailed experimental survey on WordNet. Tech. Rep. TR-2016-  
562 01, UNED, [http://e-spacio.uned.es/fez/view/bibliuned:](http://e-spacio.uned.es/fez/view/bibliuned:DptoLSI-ETSI-Infomes-Jlastra-refinement)  
563 [DptoLSI-ETSI-Infomes-Jlastra-refinement](http://e-spacio.uned.es/fez/view/bibliuned:DptoLSI-ETSI-Infomes-Jlastra-refinement).
- 564 [38] Lastra-Díaz, J. J., García-Serrano, A., Nov. 2016. HESML\_vs\_SML: scal-  
565 ability and performance benchmarks between the HESML V1R2 and  
566 SML 0.9 semantic measures libraries. Mendeley Data, v1, [http://doi.](http://doi.org/10.17632/5hg3z85wf4.1)  
567 [org/10.17632/5hg3z85wf4.1](http://doi.org/10.17632/5hg3z85wf4.1).
- 568 [39] Lastra Díaz, J. J., García Serrano, A., Jun. 2016. System and method  
569 for the indexing and retrieval of semantically annotated data using an  
570 ontology-based information retrieval model. United States Patent and  
571 Trademark Office (USPTO) application US2016/0179945 A1.
- 572 [40] Lastra-Díaz, J. J., García-Serrano, A., 2016. WNSimRep: a framework and  
573 replication dataset for ontology-based semantic similarity measures and  
574 information content models. Mendeley Data v1, [http://doi.org/10.](http://doi.org/10.17632/mpr2m8pycs.1)  
575 [17632/mpr2m8pycs.1](http://doi.org/10.17632/mpr2m8pycs.1).
- 576 [41] Lastra-Díaz, J. J., García-Serrano, A., 2016. WordNet-based word similar-  
577 ity reproducible experiments based on HESML V1R1 and ReproZip.  
578 Mendeley Data, v1, <http://doi.org/10.17632/65pxgskhz9.1>.
- 579 [42] Lastra-Díaz, J. J., García Serrano, A., 2018. HESML V1R4 Java software  
580 library of ontology-based semantic similarity measures and information  
581 content models. Mendeley Data, v4, [http://dx.doi.org/10.17632/](http://dx.doi.org/10.17632/t87s78dg78.4)  
582 [t87s78dg78.4](http://dx.doi.org/10.17632/t87s78dg78.4).
- 583 [43] Lastra-Díaz, J. J., García-Serrano, A., Batet, M., Fernández, M., Chiri-  
584 gati, F., Jun. 2017. HESML: a scalable ontology-based semantic sim-  
585 ilarity measures library with a set of reproducible experiments and a  
586 replication dataset. Information Systems 66, 97–118.
- 587 [44] Lastra-Díaz, J. J., Goikoetxea, J., Hadj Taieb, M. A., García-Serrano, A.,  
588 Ben Aouicha, M., Agirre, E., Oct. 2019. A reproducible survey on word  
589 embeddings and ontology-based methods for word similarity: linear  
590 combinations outperform the state of the art. Engineering Applications  
591 of Artificial Intelligence 85, 645–665.

- 592 [45] Lastra-Díaz, J. J., Goikoetxea, J., Hadj Taieb, M. A., García-Serrano, A.,  
593 Ben Aouicha, M., Agirre, E., Aug. 2019. Reproducibility dataset for  
594 a large experimental survey on word embeddings and ontology-based  
595 methods for word similarity. *Data in Brief*.
- 596 [46] Lastra-Díaz, J. J., Goikoetxea, J., Hadj Taieb, M. A., García-Serrano, A.,  
597 Ben Aouicha, M., Agirre, E., 2019. Word similarity benchmarks of re-  
598 cent word embedding models and ontology-based semantic similarity  
599 measures. *e-cienciaDatos*, v1, <http://dx.doi.org/10.21950/AQ1CVX>.
- 600 [47] Luong, T., Socher, R., Manning, C. D., 2013. Better word representations  
601 with recursive neural networks for morphology. In: *Proc. of CoNLL*.  
602 pp. 104–113.
- 603 [48] Merkel, D., Mar. 2014. Docker: Lightweight Linux Containers for Consistent  
604 Development and Deployment. *Linux Journal 2014 (239)*, Article  
605 No. 2.
- 606 [49] Mikolov, T., Chen, K., Corrado, G., Dean, J., May 2013. Efficient Estima-  
607 tion of Word Representations in Vector Space. arXiv:1301.3781.
- 608 [50] Miller, G. A., 1995. WordNet: A Lexical Database for English. *Communi-  
609 cations of the ACM 38 (11)*, 39–41.
- 610 [51] Miller, G. A., Charles, W. G., 1991. Contextual correlates of semantic  
611 similarity. *Language and cognitive processes 6 (1)*, 1–28.
- 612 [52] Mrkšić, N., Vulić, I., Séaghdha, D. Ó., Leviant, I., Reichart, R., Gašić,  
613 M., Korhonen, A., Young, S., 2017. Semantic Specialisation of Dis-  
614 tributional Word Vector Spaces using Monolingual and Cross-Lingual  
615 Constraints. *Trans. of the ACL 5*, 309–324.
- 616 [53] Munafò, M. R., Nosek, B. A., Bishop, D. V. M., Button, K. S., Chambers,  
617 C. D., du Sert, N. P., Simonsohn, U., Wagenmakers, E.-J., Ware, J. J.,  
618 Ioannidis, J. P. A., Jan. 2017. A manifesto for reproducible science.  
619 *Nature Human Behaviour 1*, 0021.
- 620 [54] Pedersen, T., 2008. Empiricism Is Not a Matter of Faith. *Computational  
621 Linguistics 34 (3)*, 465–470.
- 622 [55] Pedersen, T., 2008. WordNet-InfoContent-3.0.tar dataset reposi-  
623 tory. [https://www.researchgate.net/publication/273885902\\_](https://www.researchgate.net/publication/273885902_WordNet-InfoContent-3.0.tar)  
624 [WordNet-InfoContent-3.0.tar](https://www.researchgate.net/publication/273885902_WordNet-InfoContent-3.0.tar).
- 625 [56] Pedersen, T., Nov. 2013. Measuring the Similarity and Relatedness of Con-  
626 cepts: a MICAI 2013 Tutorial.
- 627 [57] Pennington, J., Socher, R., Manning, C. D., 2014. Glove: Global vectors  
628 for word representation. *Proc. of EMNLP 12*, 1532–1543.
- 629 [58] Pirró, G., Nov. 2009. A semantic similarity metric combining features and  
630 intrinsic information content. *Data & Knowledge Engineering 68 (11)*,  
631 1289–1308.

- 632 [59] Radinsky, K., Agichtein, E., Gabrilovich, E., Markovitch, S., Mar. 2011. A  
633 word at a time: computing word relatedness using temporal semantic  
634 analysis. In: Proc. of the Intl. Conf. on WWW. ACM, pp. 337–346.
- 635 [60] Rubenstein, H., Goodenough, J. B., Oct. 1965. Contextual Correlates of  
636 Synonymy. *Communications of the ACM* 8 (10), 627–633.
- 637 [61] Schwartz, R., Reichart, R., Rappoport, A., 2015. Symmetric pattern based  
638 word embeddings for improved word similarity prediction. In: Proc. of  
639 the Conf. on Computational Natural Language Learning. pp. 258–267.
- 640 [62] Sefraoui, O., Aissaoui, M., Eleuldj, M., 2012. OpenStack: toward an open-  
641 source solution for cloud computing. *International Journal of Computer*  
642 *Applications in Technology* 55 (3), 38–42.
- 643 [63] Szumlanski, S. R., Gomez, F., Sims, V. K., Aug. 2013. A New Set of  
644 Norms for Semantic Relatedness Measures. In: Proc. of the 51st Annual  
645 Meeting of the Association for Computational Linguistics (ACL'2013).  
646 Vol. 2. [aclweb.org](http://aclweb.org), Sofia, Bulgaria, pp. 890–895.
- 647 [64] Wieling, M., Rawee, J., van Noord, G., Dec. 2018. Reproducibility in Com-  
648 putational Linguistics: Are We Willing to Share? *Computational Lin-*  
649 *guistics* 44 (4), 641–649.
- 650 [65] Wieting, J., Bansal, M., Gimpel, K., Livescu, K., Roth, D., 2015. From  
651 Paraphrase Database to Compositional Paraphrase Model and Back.  
652 *Trans. of the ACL* 3, 345–358.
- 653 [66] Wolke, A., Bichler, M., Chirigati, F., Steeves, V., Jul. 2016. Reproducible  
654 experiments on dynamic resource allocation in cloud data centers. *In-*  
655 *formation Systems* 59, 98–101.
- 656 [67] Yang, D., Powers, D. M., 2006. Verb similarity on the taxonomy of Word-  
657 Net. In: Proc. of the 3th Intl. WordNet Conf. (GWC). pp. 121–128.
- 658 [68] Yeh, E., Ramage, D., Manning, C. D., Agirre, E., Soroa, A., 2009. Wiki-  
659 Walk: Random Walks on Wikipedia for Semantic Relatedness. In: Proc.  
660 of the 2009 Workshop on Graph-based Methods for Natural Language  
661 Processing. TextGraphs-4. ACL, Stroudsburg, PA, USA, pp. 41–49.