Learning Disentangled Representations in Natural Language Definitions with Semantic Role Labeling Supervision

Anonymous ACL submission

Abstract

Disentangling the encodings of neural models is a fundamental aspect for improving interpretability, semantic control and downstream task performance in Natural Language Processing. Currently, most disentanglement methods are unsupervised or rely on synthetic datasets with known generative factors. We argue that 800 recurrent syntactic and semantic regularities in textual data can be used to provide the models with both structural biases and generative factors. We leverage the semantic structures present in a representative and semantically dense category of sentence types, definitional sentences, for training a Variational 015 Autoencoder to learn disentangled representations. Our experimental results show that the proposed model outperforms unsupervised baselines on several qualitative and quantitative benchmarks for disentanglement, and it also improves the results in the downstream task of definition modeling.

1 Introduction

004

011

017

034

040

Learning disentangled representations is a fundamental step towards enhancing the interpretability of the encodings in deep generative models, as well as improving their downstream performance and generalization ability. Disentangled representations aim to encode the fundamental structure of the data in a more explicit manner, where independent latent variables are embedded for each generative factor (Bengio et al., 2013).

Previous work in machine learning proposed to learn disentangled representations by modifying the ELBO objective of the Variational Autoencoders (VAE) (Kingma and Welling, 2014), within an unsupervised framework (Higgins et al., 2017; Kim and Mnih, 2018; Chen et al., 2018). On the other hand, a more recent line of work claims the benefits of supervision in disentanglement (Locatello et al., 2019) and it advocates the importance of designing frameworks able to exploit structures



Evaluation



Figure 1: Left: Supervision mechanism with definition semantic roles (DSR) encoded in the latent space. The dotted arrow represent the conditional VAE version. Right: Evaluation framework.

in the data for introducing inductive biases. In parallel, disentanglement approaches for NLP have been tackling text style transfer, and evaluating the results with extrinsic metrics, such as style transfer accuracy (Hu et al., 2017; John et al., 2019; Cheng et al., 2020).

While style transfer approaches investigate the ability to disentangle and control syntactical factors such as tense and gender, the aspect of understanding and disentangling the semantic structure in language is under-explored. Furthermore, evaluating disentanglement is challenging, because it requires knowledge of generative factors, leading most approaches to train on synthetic datasets (Higgins et al., 2017; Zhang et al., 2021).

In this work, we argue that recurrent semantic structures at sentence level can be leveraged both as inductive biases for enhancing disentanglement (RQ1) but also for providing meaningful generative factors that can be employed for evaluating the degree of disentanglement (RQ2). We also investigate whether organizing the generative factors in groups may facilitate learning and disentanglement (RQ3). As a result, this work focuses on natural

107

108

110

111

112

113

language definitions, which are a textual resource characterised by a principled structure in terms of semantic roles, as demonstrated by previous work which proposed the extraction of structural and semantic patterns in this kind of data (Silva et al., 2016, 2018).

Seeking to address the highlighted issues and answer the research questions, we make the following contributions, also depicted in Figure 1.

1) We design a supervised framework for enhancing disentanglement in language representations by conditioning on the information provided by the semantic role labels (SRL) in natural language definitions. We present two mechanisms for injecting SRL biases into latent variables, firstly, reconstructing both words and corresponding SRL in a VAE, secondly, employing SRL information as input variable for a Conditional VAE (Zhao et al., 2017).

2) We propose the first framework for evaluating the disentanglement properties of the encodings on non-synthetic textual datasets. Our evaluation framework employs semantic role labels as generative factors, enabling the measurement of several contemporary quantitative metrics. The results show that the proposed bias injection mechanisms are able to increase the degree of disentanglement of the representations.

3) We demonstrate that models trained with our disentanglement framework are able to outperform contemporary baselines in the downstream task of definition modeling (Noraset et al., 2017).

Disentangling framework 2

In this section we first describe the framework that we designed for improving disentanglement in natural language definitions with semantic role labels. Secondly, we present three models, shown in Figure 2 based on the Variational Autoencoder (VAE) (Bowman et al., 2016) for achieving disentanglement.

Disentangling definitions 2.1

Definition semantic roles Our framework is based on natural language definitions, which are a particular type of linguistic expression, characterised by high abstraction, and specific phrasal properties. Previous work in NLP for dictionary definitions (Silva et al., 2018) has shown that there are categories that can be consistently found in most definitions. In fact, Silva et al. (2018) define 114

precise Semantic Role Labels (SRL) for phrases representing definitions, under the name of Definition Semantic Roles (DSR).

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

163

The example from (Silva et al., 2018) classifies the semantic roles within "english poets who lived in the lake district" as follows. "poets" as noun category (supertype), "english" as quality of the term (Differentia Quality), "who lived" as event that the subject is involved with (differentia event), and "in the lake district" as the location of the action (Event location). The full DSRs proposed by Silva et al. (2018) are reported in Table 7 in Appendix A.

Disentangling using SRL Our goal is to enhance disentanglement in natural language by injecting categorical structures into latent variables. We find that this goal is well aligned with the findings of Locatello et al. (2019), where it is claimed that a higher degree of disentanglement may benefit from supervision and inductive biases. Our hypothesis is that we may leverage such semantic information for learning representation with higher degree of disentanglement. While in the context of this work we use dictionary definitions as a target empirical setting, we conjecture that these conclusions can be extended to broader definitional sentence-types.

2.2 **Definition VAEs**

Unsupervised VAE The first baseline model that we consider is the traditional variational autoencoder (VAE) for sentences (Bowman et al., 2016), which operates in an unsupervised fashion, as in Figure 2a. The unsupervised VAE employs a multivariate gaussian prior distribution p(z) and generates a sentence x with a decoder network $p_{\theta}(x|z)$. The joint distribution for the decoder is defined as $p(z)p_{\theta}(x|z)$, which, for a sequence of tokens x of length T result as $p_{\theta}(x|z) = \prod_{i=1}^{T} p_{\theta}(x_i|x_{< i}, z)$. The VAE objective consists into maximizing the expectation of the log-likelihood which is defined as $\mathbb{E}_{p(x)} \log p_{\theta}(x)$. Due to the computational intractability of the such expectation value, the variational distribution q_{θ} is employed to approximate $p_{\theta}(z|x).$

As a result, an evidence lower bound \mathcal{L}_{VAE} (ELBO) where $\mathbb{E}_{p(x)}[\log p_{\theta}(x)] \geq \mathcal{L}_{VAE}$, is derived as follows:

$$\mathcal{L}_{\text{Tokens}} = \mathbb{E}_{q_{\phi}(z|x)} \Big[\log p_{\theta}(x|z) \Big]$$
(1)
- KL $q_{\phi}(z|x) || p(z)$

DSR supervised VAE The aim of this model is to inject the categorical structure of the definition



Figure 2: Proposed architectures for learning disentangled representations in definitions.

semantic roles (DSR) into the latent variables, by 164 factorizing them into the VAE auto-encoding ob-165 jective function. In order to achieve this goal, we 166 introduce the variable r for semantic roles, and train 167 the "DSR VAE", where both sentence and semantic 168 169 roles are auto-encoded. As a result, two separate losses are produced and added together for the final loss, as shown in Figure 2b. The ELBO for 171 semantic roles is defined as follows: 172

$$\mathcal{L}_{\text{Roles}} = \mathbb{E}_{q_{\phi}(z|r)} \Big[\log p_{\theta}(r|z) \Big]$$

$$- \mathrm{KL}q_{\phi}(z|r) || p(z)$$
(2)

173

174

175

176

177

178

179

180

181

185

186

189

190

191

192

194

195

196

198

199

200

The final loss is given by $\mathcal{L}_{Tokens} + \mathcal{L}_{Roles}$.

Conditional VAE with SRL For explicitly leveraging the definition semantic roles, we propose a supervision mechanism based on the Conditional VAE (CVAE) (Zhao et al., 2017), shown in Figure 2c. Similar to the previously described model, we instantiate a VAE framework, where x is the variable for the tokens, and r for the roles. We perform auto-encoding for both roles and tokens, and additionally, we condition the decoder network on the roles. The CVAE is trained to maximize the conditional log likelihood of x given r, which involves an intractable marginalization over the latent variable z.

The ELBO is defined as:

$$\mathcal{L}_{\text{CVAE}} = \mathbb{E}_{q_{\phi}(z|r,x)} \Big[\log p_{\theta}(x|z,r) \Big]$$
(3)
- KL $q_{\phi}(z|x,r) || p(z|r)$

Training The training process follows the variational autoencoding methodology (Kingma and Welling, 2014). First, tokenization is performed in the sentences and the roles. The Encoder network involves feeding both first into embedding layers, then into LSTM layers. Subsequently, two vectors μ and σ are sampled with two linear layers, and the vector z is computed with the re-parameterization trick. Finally, the decoder network is built with LSTM and another embedding layer, which return the same dimension that was given as input.

3 Evaluation framework

We first present the evaluation framework that for measuring disentanglement, then describe and justify the generative factor setup used in the experiments.

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

221

224

225

226

227

229

230

231

232

233

234

237

238

239

240

241

242

243

3.1 DSR as generative factors

While early approaches for disentanglement in NLP have been proposed in the context of in style transfer applications (John et al., 2019; Cheng et al., 2020) and are assessed purely in terms of style transfer accuracy, evaluating the intrinsic properties of the latent encodings is fundamental for disentanglement, as mentioned in several machine learning approaches (Higgins et al., 2017; Kim and Mnih, 2018). Recently, Zhang et al. (2021) proposed a framework for computing several popular quantitative disentanglement metrics such as (Higgins et al., 2017; Kim and Mnih, 2017; Kim and Mnih, 2018) testing it on synthetic datasets. The limitation in (Zhang et al., 2021) is that is works only with synthetic datasets.

In this work, we propose a method where semantic role labels, such as the ones provided in (Silva et al., 2018), are used as generative factors for evaluating the degree of disentanglement in the encodings. The framework, overview in Figure 3, considers multiple generative factors, where each factor is composed by a number of semantic roles (for example the factor "location" includes, origin-location, and event-location). In this way, the dataset can be seen as the result of a sampling of multiple generative factors, which is the same principle used when creating synthetic datasets for disentanglement. Once the generative factors are defined, the framework is enabled to compute a number of quantitative metrics for disentanglement, following the work from Zhang et al. (2021).

3.2 Semantics and Syntax groups of DSR

In order to categorize the definition semantic roles (DSR), we consider its structural and semantic dimensions in terms of their contribution to either the meaning (e.g., quality, location) or the structure

Group 1: Semantics	Group 2: Syntax	Group 3: Semantics	Group 4: Syntax
Supertype: SUPERTYPE Quality: DIFFERENTIA-QUALITY Event: DIFFERENTIA-EVENT	Supertype: SUPERTYPE Main	Quality DIFFERENTIA-QUALITY QUALITY-MODIFIER ACCESSORY OUTALITY	Main DIFFERENTIA-QUALITY DIFFERENTIA-EVENT
Location EVENT-LOCATION ORIGIN-LOCATION	DIFFERENTIA-EVENT Modifier Event	Event DIFFERENTIA-EVENT EVENT_TIME	Modifier Event EVENT-LOCATION EVENT-TIME
Modifier QUALITY-MODIFIER EVENT-TIME	EVENT-TIME ASSOCIATED-FACT	Location EVENT-LOCATION	Modifier Quality QUALITY-MODIFIER PURPOSE
Statement PURPOSE ASSOCIATED-FACT	QUALITY-MODIFIER PURPOSE	Statement PURPOSE	ASSOCIATED-FACT Accessory ACCESSORY-DETERMINER
Accessory ACCESSORY-DETERMINER ACCESSORY-QUALITY	Accessory ACCESSORY-DETERMINER ACCESSORY-QUALITY	ASSOCIATED-FACT ACCESSORY-DETERMINER	ACCESSORY-QUALITY

Figure 3: Generative factors for definitions.

244 (e.g., main terms, modifiers) of the definition sentence. We first create two DSR groups with seman-245 tic and two based on syntax, to evaluate which one 246 would better facilitate disentanglement. For both syntax and semantic, we then create a groups with 248 "supertype" DSR and one without it, in order to 249 understand the impact of the supertype DSR. The importance of "supertype" is due to its contribu-251 tion to both abstraction groups and its predominant presence on the datasets analyzed ($\geq 97\%$).

Group 1: Semantics with Supertype Sets the factors in terms of their meaning, essentially abstracting categories of the DSRs, including the SU-PERTYPE DSR as a single factor. Qualification, location, modification, declaration (statement) and supplementation (accessory) are semantic roles of a given term to its definition, which are described by the DSRs. For example, "event location" and "origin location" are inserted into the Location factor, while "purpose" and "associated fact" are in statement.

257

258

260

261

262

263

267

271

273

Group 2: Syntax with Supertype Sets the fac-265 tors in terms of their structural role in the definition 266 sentence, including the SUPERTYPE DSR as a single factor. The ORIGIN-LOCATION DSR is 268 omitted due to its syntactic overlap with EVENT-LOCATION and its low frequency in the datasets. Typically, a "main" term is followed by a set of modifiers (event, quality), which may be comple-272 mented by accessory terms.

Group 3: Semantics without Supertype Similar to group 1, but excluding the SUPERTYPE DSR, and repositioning the factor from *modifier* and accessory for higher abstraction. The rela-277 tions of modification and supplementation (present in group 1) are suppressed to focus on lexi-279 cal semantics, moving the label ACCESSORY- DETERMINER to the declaratory group, EVENT-TIME to the event group and all quality related labels to the qualification group.

281

283

284

287

289

290

291

292

293

294

295

296

297

298

300

301

302

303

304

305

307

308

309

310

311

312

313

314

315

316

Group 4: Syntax without Supertype Similar to group 2, but excluding the SUPERTYPE DSR. Further abstractions are not conducted, as the definition roles already offer a stable structure for sentence construction.

4 **Related work**

Disentangled VAEs in language Early approaches in text disentanglement use VAEs with multiple adversarial losses for style transfer (Hu et al., 2017; John et al., 2019). More recently, Cheng et al. (2020) propose a style transfer method which minimizing the mutual information between the latent and the observed variable, while Colombo et al. (2021) propose an upper bound of mutual information for fair text classification. On the other hand, we propose to disentangle representation using biases provided as semantic roles and design two VAE models to inject structural semantic information into the representation.

Disentanglement Evaluation Vishnubhotla et al. (2021) evaluate disentanglement in synthetic text on various NLP tasks such as classification, retrieval and style transfer. Zhang et al. (2021) evaluate disentanglement of various VAE models on synthetic datasets where generative factors are known. Differently from these methods, we propose a new framework to evaluate non-synthetic natural languages, where semantic role labels are used as generative factors. We model linguistic features of natural language definitions, with the goal of exploring the semantic properties that are encapsulated in it.

Definition models Early approaches in definition encoding include (Hill et al., 2016), which pro-

pose the first neural embedding model for dictionar-317 ies, and (Bahdanau et al., 2017), which present an 318 RNN-based encoder decoder architecture for tex-319 tual entailment and reading comprehension. More recently, methods based on Autoencoders (Bosc 321 and Vincent, 2018) and transformers (Tsukagoshi 322 et al., 2021) have been proposed. Various ap-323 proaches for the task of generating a definition from a word (Definition Modeling) have been proposed, including RNN-based methods (Noraset et al., 2017), soft attention mechanisms (Gadetsky et al., 2018), and span-based encoding schemes (Bevilac-328 qua et al., 2020). The semantic aspect of natural 329 language definitions are explored in (Silva et al., 2016, 2018), where the concept of definition se-331 mantic roles is proposed.

5 Empirical analysis

333

336

337

339

340

341

342

347

351

352

353

355

356

357

361 362

363

In this section, we firstly describe the empirical setup for experiments, secondly, we provide qualitative evaluation and thirdly, we measure various quantitative metrics. Finally, we demonstrate the capacity of the proposed models in the downstream task of definition modeling. The full experimental pipeline is available in a code package submitted as supplementary material (Appendix D).

5.1 Experimental setup

Datasets Definition sentences and their respective semantic role structures are sourced from three different datasets by (Silva et al., 2016) with the characteristics described in Table 1. All datasets are automatically annotated with DSR tags for each token, using the method proposed by (Silva et al., 2016). The datasets differ not only in sentence length and size, but also in textual style: while WordNet and Wiktionary sentences tend to be formatted as dictionary definitions, Wikipedia sentences are lengthier and less adherent to a typical definition structure.

Hyperparameter choices Experiments are conducted to cover a set of 3 hyperparameters: First, the VAE architecture used: 1) Unsupervised VAE
2) Supervised with SRL 3) CVAE with SRL. Second, the generative factor grouping, which includes:
1) Semantic w/ supertype 2) Syntactic w/ supertype
3) Semantic w/o supertype 4) Syntactic w/o supertype. Third, the dimensionality of VAE latent representation (z): 4, 5, 7, 128.

The choice of architecture allows evaluation of the impact of DSR label conditioning in two dis-

Dataset	Num sents.	Avg. length	Version
Wordnet	93,699	9	WordNet 3.0
Wiktionary	464,243	8	Dec, 2016
Wikipedia	1,500,323	12	Dec, 2016

Table 1: Statistics from definition datasets.

tinct ways: as part of the autoencoding objective function, and as a conditional variable of the decoder, addressing our research questions **RQ1** and **RQ2**. The choice of generative factor grouping can indicate the best ways to organize the factors, addressing **RQ3**.

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

The dimensionality of the representation is set to match the number of generative factors, in an attempt to force disentanglement by alignment of each dimension to a single factor. The dimension sizes are then defined to be 4 (alignment with groupings 3 and 4), 5 (alignment with grouping 2) or 7 (alignment with grouping 1). However, different levels of disentanglement can be achieved with mismatching dimensions and factors. So all possible combinations of factors and representation sizes are tested and a size of 128 is included to evaluate the impact of a higher number of parameters in each grouping.

Implementation Details Neural Network hyperparameters are kept fixed for all quantitative experiments, with the following values, based on a previous experiment from (Shen et al., 2020). (1) Number of hidden layers: 1, (2) Dimension of the hidden layer: 512, (3) VAE $\lambda_{KL} = 0.1$, (4) Epochs=20, (5) Batch size=32 for Wikipedia, 64 for the rest. All VAEs built for the experiments are composed of a LSTM as sequence encoder, a hidden layer, an embedding (representation) layer, and a LSTM decoder. Dropout (20%) is done for both encoder and decoder inputs. To provide the inputs and outputs for the VAEs, the definition sentences are tokenized into sub-words with a Byte Pair Encoding (BPE) scheme, and converted into token embeddings with the T5 transformer model (Raffel et al., 2020), with an embedding size of 512. The use of transformer embeddings introduce richer semantic information that can be leveraged by the VAE in the construction of its representations.

5.2 Qualitative Evaluation

We evaluate the representations of the trained models in terms of their disentanglement, by analysing 1) traversals of the latent space, 2) encoding interpolation 3) encoding visualization with dimension-

VAE	a surgical procedure for one purpose a parasitic procedure for one skull a parasitic procedure for its content
DSR	a fictional character having close incense a fictional name consisting by the kitchen a fictional name consisting of the brothers
CVAE	a simple scheme where members also must take minerals a simple scheme where members also must be out in Renaissance a simple scheme where members also specialized on a tract

Table 2: Traversal showing disentangled and entangled factors in Wordnet.

410	ality	reduction.

411

413

415

417

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

Latent space traversals Traversal evaluation is a standard procedure with image disentangle-412 ment (Higgins et al., 2017; Kim and Mnih, 2018). The traversal of a latent factor is obtained as the 414 decoding of the vectors corresponding to the latent variables, where the evaluated factor is changed 416 within a fixed interval, while all others are kept fixed. If the representation is disentangled, when 418 a latent factor is traversed, the decoded sentences 419 should only change with respect to that factor. This 420 means that after training the model we are able to probe the representation for each latent variable.

> In Figure 2 we report examples from Wordnet obtained for 128 latent variables for VAE and DSR VAE and CVAE. We observe that all models show some degree of disentanglement, because the decoded sentences only change few attributes, denoting control. In particular, the CVAE is the one with the more disentangled representations, given that it maintains the first phrase constant, and varies the second one smoothly.

Interpolation In this experiment, we demonstrate the ability of autoencoder models to provide smooth transition between latent space representations of sentences (Bowman et al., 2016). In practice, the interpolation mechanism takes two sentences x_1 and x_2 , and uses their posterior mean as the latent features z_1 and z_2 , respectively. It interpolates a path $z_t = z_1 \cdot (1-t) + z_2 \cdot t$ with t increased from 0 to 1 by a step size of 0.1. As a result, 9 sentences are generated on each interpolation step.

In Table 3 we provide qualitative results with latent space interpolation on Wordnet. Interestingly, the DSR-supervised VAE shows ability to paraphrase semantic concepts, for example bridging the concept of teaching and learning on the starting sentence with the concept of train and loading goods we find in the middle the notion



Table 3: Interpolation examples in Wordnet.



Figure 4: t-SNE plot of 9370 definition representations (128 dimensions), generated from Unsupervised VAE (U), DSR supervision (S) and Conditional VAE (C).

of "organising information" and "control components" which are noteworthy semantic bridges. This type of localised semantic control provided by the operations of traversal and interpolation over intensional-level (definitional) sentences can potentially support quasi-symbolic operations over the latent space.

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

t-SNE plot t-SNE (t-distributed Stochastic Neighbor Embedding) (Van der Maaten and Hinton, 2008) is a popular method for non-linear dimensionality reduction, that allows the visualization of complex high-dimensional feature spaces, such as the representation space produced by a VAE. Figure 4 presents a 2D plot of t-SNE transformations for each one of the evaluated models, from which the clustering of DSR patterns can be observed. While the supervision with DSR labels promotes clustering of the patterns around the center of the plot, cVAE compacts the cluster on the edges, allowing better separation of the DIFF-QUALITY+DIFF-EVENT and DIFF-EVENT+DIFF-QUALITY patterns. From the cVAE plot is also possible to visualise a smoother transition between the two major patterns observed



Figure 5: Metrics mean grouped.

in the dataset: from DIFFERENTIA-QUALITY
(red) to DIFFERENTIA-EVENT (blue), with the
combination of both patterns being coloured purple/violet. UMAP transformations are also performed and the plots are presented in the supplemental material (Appendix C).

5.3 Quantitative Evaluation

In this experiment we probe the representation learned by the proposed VAE models using eight popular quantitative metrics for disentanglement, namely: z-diff (Higgins et al., 2017), z-minvar (Kim and Mnih, 2018), Mutual Information Gap (MIG) (Chen et al., 2018), Modularity & Explicitness (Ridgeway and Mozer, 2018), and from (Eastwood and Williams, 2018)(disentanglement, completeness, informativeness). Further details about the metrics are provided in Appendix B.

Experimental Setup We evaluate VAE (U), DSR 491 VAE (S) and CVAE (C) on Wordnet (WN), Wik-492 tionary (WT) and Wikipedia (WP) datasets. Evalu-493 ation is performed under the framework explained 494 in Section 3. Each combination of VAE architec-495 496 ture, generative factor grouping and representation size was trained and quantitatively tested, by cal-497 culating the previously mentioned disentanglement 498 metrics. For computing the metrics we follow the 499 experiments of Zhang et al. (2021). 500

D	z-diff		z-min-var ↓		MIG			Modularity				
	U	S	С	U	S	С	U	S	С	U	S	С
WN	70.0	69.1	77.0	48.2	50.3	53.2	.067	.057	.059	.793	.804	.765
WT	59.7	61.9	63.5	40.0	38.5	43.0	.112	.095	.065	.535	.424	.629
WP	57.5	63.0	64.7	39.8	38.6	42.0	.046	.041	.037	.771	.745	.757
D	Ex	plicitn	ess	Diser	ntangle	ement	Con	nplete	ness	Infor	mative	ness ↓
	U	S	С	U	S	С	U	S	С	U	S	С
WN	.519	.532	.527	.022	.021	.031	.013	.013	.017	.364	.361	.399
WT	.584	.593	.616	.014	.011	.013	.013	.013	.011	.377	.373	.385
WP	.545	.557	.600	.007	.007	.005	.007	.007	.004	.375	.373	.374

Table 4: Quantitative disentanglement metrics.

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

Analysis The results presented in Table 4 show that across all datasets, the application of DSR categories as biases results in a measurable improvement in disentanglement (RQ1) and that the use of DSRs as generative factors produces meaningful disentangled representations (RQ2). More specifically, z-diff presents the highest and most consistent improvement, specially with the CVAE, indicating higher interpretability when inferring single generative factors from the representations. Explicitness results are also consistent, indicating higher coverage of each factor. Improvements on Modularity, Disentanglement Score, Completeness and Informativeness are less consistent, indicating that the factors share substantial information between them. On the other hand, z-min-var, MIG counter the trend of improvement, due to the fact that they are designed to strongly penalize nonalignment of single pairs < factor \leftrightarrow latent dimension> (e.g., linear combinations). As a result, they penalize the existence of dependency and hierarchy relations which is present in most DSR categories, e.g., DIFFERENTIA-EVENT \rightarrow EVENT-TIME.

We analyse how semantic groupings affect disentanglement in Figure 5b (RQ3). Overall, we notice that syntax based groups have higher disentanglement, indicating that it is easier to disentangle syntactic phrase components. For Modularity the result is the opposite, indicating that semantic groupings promote higher independence between factors. Following (Zhang et al., 2021), the values in Figure 5b for the metrics Completeness and Disentanglement score are multiplied by 10, in order to facilitate the visualization.

Finally, we find that a low number of latent dimensions leads to smaller degree of disentanglement. The experiments with 4,5,7 and 128 latents are reported in Figure 5a.

5.4 Definition Generation

In this experiment, we assess the proposed VAE models in the task of "Definition Modeling" (No-

480

481

482

483

484

485

486

487

488

489

Word	Definition Model	Unsupervised VAE	Supervised VAE
repulse	the act of making a gun	the act of moving forward	act in a hostile state
colonise	make a new or vital part	the state of being in a particular place	settle or cause to be easily removed
involve	make a specific purpose	make a specific effect	a specific act of making something
mitochondrion	a cell that is used to treat the blood	a substance that is used to treat a body reaction	a cell that is a source of an organic process
heat	a change in the surface of a liquid	a sudden increase in the flow of heat	a sudden increase in the temperature

Table 5: Definition generation examples for the Wordnet dataset.

Perplexity \downarrow				Bleu		
Data	DM	VAE	DSR	DM	VAE	DSR
WN	88.59	80.36	80.27	9.12	10.27	10.26
WT	42.51	39.09	38.64	6.70	7.53	7.59
WP	13.09	12.39	12.47	11.89	12.32	12.34

Table 6: Quantitative metrics for definition generation.

raset et al., 2017), where the goal is to generate a natural language definition given the word to be defined (definiendum).

Experimental setup During training, we adopt the "seed" setup (Noraset et al., 2017), which involves providing the definiendum concatenated with the definition tokens as input for the model. At generation time, the model takes as input only the word which needs to be defined, and leverages a trained model for computing the definition latent encoding. Such encoding is then fed into a softmax function and subsequently a multinomial probability distribution is sampled for decoding the latent variable into the final definition sentence.

We compare the proposed unsupervised and DSR-supervised VAEs with the LSTM-based Definition-model approach from (Gadetsky et al., 2018), both using the "seed" setup. The CVAE is not explored in this experiment in order to have a more fair comparison with the Definition model. We train the baseline and our models with similar setups, following (Gadetsky et al., 2018). We perform language model pretraining on the WikiText-103 dataset (Merity et al., 2016) for 1 epoch, then train on the downstream dataset for 10 epochs. Additionally, all models are initialised using Google Word2Vec pretrained vectors, following (Gadetsky et al., 2018).

571**Results** We report the perplexity and Bleu (Pap-572ineni et al., 2002) results in Table 6. We observe573that the proposed variational autoencoder models574achieve an improvement on both perplexity and575Bleu compared to the RNN baseline. The DSR576VAE achieves the best perplexity and bleu on 2577out of 3 datasets while the unsupervised VAE is578the best performing model in the other cases. We

justify the success of VAE models due to their disentangling properties, and also their ability to learn smooth encodings, a benefit deriving from sampling variable for re-parameterization. In particular, we attribute the success of the DSR VAE to the additional knowledge that has been injected into its latent variables. 579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

600

601

602

603

604

605

607

608

609

610

611

612

613

614

615

616

617

618

Some generation examples from the Wordnet dataset are provided in Table 5. Such examples show that the proposed VAE models are able to leverage the structural and semantic information of the learned definition roles to better approximate the defined concept. In particular, we notice some semantically strong linguistic elements in the definitions decoded with DSR supervision, for example DSR is the only model able to link the verb "repulse" with the hostile adjective, the verb colonise with the similar verb "settle", and the word "heat" with temperature.

The strong performance in this definition generation task demonstrates that the disentangled representations have provided the VAE models with significant generalization capability, confirming that disentangling is beneficial for various applications tasks.

6 Discussion

We propose a novel VAE-based framework for learning and evaluating disentangled representations in natural language definitions. We leverage the semantic structure present in dictionaries as inductive biases for improving disentanglement in VAEs, and as generative factors during evaluation. Our evaluation shows, both with qualitative investigations and with quantitative metrics, that the proposed framework is able to produce encodings with a higher degree of disentanglement. Finally, our models outperform existing baselines on a definition modeling application, demonstrating the generalization capabilities of disentangled representations.

563

564

566

567

569

570

543

544

References

619

622

623

624

625

626

627

629

631

634

635

641

647

652

663

664

666

667

670

673

- Dzmitry Bahdanau, Tom Bosc, Stanisław Jastrzębski, Edward Grefenstette, Pascal Vincent, and Yoshua Bengio. 2017. Learning to compute word embeddings on the fly. arXiv preprint arXiv:1706.00286.
 - Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. 2013. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35:1798–1828.
- Michele Bevilacqua, Marco Maru, and Roberto Navigli. 2020. Generationary or:"how we went beyond word sense inventories and learned to gloss". In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7207–7221.
- Tom Bosc and Pascal Vincent. 2018. Auto-encoding dictionary definitions into consistent word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1522–1532.
- Samuel Bowman, Luke Vilnis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and Samy Bengio. 2016. Generating sentences from a continuous space. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 10–21.
- Ricky TQ Chen, Xuechen Li, Roger Grosse, and David Duvenaud. 2018. Isolating sources of disentanglement in vaes. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 2615–2625.
- Pengyu Cheng, Martin Renqiang Min, Dinghan Shen, Christopher Malon, Yizhe Zhang, Yitong Li, and Lawrence Carin. 2020. Improving disentangled text representation learning with information-theoretic guidance. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7530–7541.
- Pierre Colombo, Pablo Piantanida, and Chloé Clavel. 2021. A novel estimator of mutual information for learning to disentangle textual representations. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics, pages 6539– 6550.
- Cian Eastwood and Christopher KI Williams. 2018. A framework for the quantitative evaluation of disentangled representations. In *6th International Conference on Learning Representations*.
- Artyom Gadetsky, Ilya Yakubovskiy, and Dmitry Vetrov. 2018. Conditional generators of words definitions. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 266–271.
- Irina Higgins, Loïc Matthey, Arka Pal, Christopher P. Burgess, Xavier Glorot, Matthew M. Botvinick, Shakir Mohamed, and Alexander Lerchner. 2017.

beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*. 674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

- Felix Hill, KyungHyun Cho, Anna Korhonen, and Yoshua Bengio. 2016. Learning to understand phrases by embedding the dictionary. *Transactions of the Association for Computational Linguistics*, 4:17– 30.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P Xing. 2017. Toward controlled generation of text. In *International Conference on Machine Learning*, pages 1587–1596. PMLR.
- Vineet John, Lili Mou, Hareesh Bahuleyan, and Olga Vechtomova. 2019. Disentangled representation learning for non-parallel text style transfer. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 424–434.
- Hyunjik Kim and Andriy Mnih. 2018. Disentangling by factorising. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2649–2658. PMLR.
- Diederik P. Kingma and Max Welling. 2014. Autoencoding variational bayes.
- Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. 2019. Challenging common assumptions in the unsupervised learning of disentangled representations. In *international conference on machine learning*, pages 4114–4124. PMLR.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- Thanapon Noraset, Chen Liang, Larry Birnbaum, and Doug Downey. 2017. Definition modeling: Learning to define word embeddings in natural language. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th annual meeting of the Association for Computational Linguistics, pages 311–318.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Karl Ridgeway and Michael C Mozer. 2018. Learning deep disentangled embeddings with the f-statistic loss. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 185–194.

728

- 736 737 740 741 742 743 744 745 746 747 748 749 750
- 751 754 755

757

- Tianxiao Shen, Jonas Mueller, Regina Barzilay, and Tommi Jaakkola. 2020. Educating text autoencoders: Latent representation guidance via denoising. In International Conference on Machine Learning, pages 8719-8729. PMLR.
 - V. S. Silva, S. Handschuh, and A. Freitas. 2016. Categorization of semantic roles for dictionary definitions. In CogALex@COLING.
 - Vivian S Silva, Siegfried Handschuh, and André Freitas. 2018. Recognizing and justifying text entailment through distributional navigation on definition graphs. In Thirty-Second AAAI Conference on Artificial Intelligence.
 - Hayato Tsukagoshi, Ryohei Sasano, and Koichi Takeda. 2021. Defsent: Sentence embeddings using definition sentences. In ACL/IJCNLP.
 - Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. Journal of machine learning research, 9(11).
 - Krishnapriya Vishnubhotla, Graeme Hirst, and Frank Rudzicz. 2021. An evaluation of disentangled representation learning for texts. In Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pages 1939-1951.
- Lan Zhang, Victor Prokhorov, and Ehsan Shareghi. 2021. Unsupervised representation disentanglement of text: An evaluation on synthetic datasets. In Proceedings of the 6th Workshop on Representation Learning for NLP (RepL4NLP-2021).
- Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017. Learning discourse-level diversity for neural dialog models using conditional variational autoencoders. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 654-664.

A Definition Semantic Roles

763

768

770

772

773

775

776

779

781

784

785

787

765

The datasets used in our experiments are introduced in (Silva et al., 2018). We report in Table 7 the annotated categories.

Role	Description
Supertype	the immediate or ancestral entity's superclass
Differentia	a quality that distinguishes the entity from the
quality	others under the same supertype
Differentia	an event (action, state or process) in which the
event	entity participates and that is mandatory to dis-
	tinguish it from the others under the same super-
	type
Event	the location of a differentia event
location	
Event time	the time in which a differentia event happens
Origin	the entity's location of origin
location	
Quality	degree, frequency or manner modifiers that con-
modifier	strain a differentia quality
Purpose	the main goal of the entity's existence or occur-
	rence
Associated	a fact whose occurrence is/was linked to the
fact	entity's existence or occurrence
Accessory	a determiner expression that doesn't constrain
determiner	the supertype / differentia scope
Accessory	a quality that is not essential to characterize the
quality	entity
Role	a particle, such as a phrasal verb complement,
particle	non-contiguous to the other role components

Table 7: Semantic Role Labels for dictionary defini-tions.

B Disentanglement Metrics

- 1. z_{diff} accuracy (Higgins et al., 2017): The accuracy of a predictor for $p(y|z_{diff}^b)$, where z_{diff}^b is the absolute linear difference between the inferred latent representations for a batch *B* of latent vectors, written as a percentage value. Higher values imply better disentanglement.
- 2. $z_{min_var} \ error$ (Kim and Mnih, 2018): For a chosen factor k, data is generated with this factor fixed but all other factors varying randomly; their representations are obtained, with each dimension normalised by its empirical standard deviation over the full data (or a large enough random subset); the empirical variance is taken for each dimension of these normalised representations. Then the index of the dimension with the lowest variance and the target index k provide one training input/output example for the classifier. Thus, if the representation is perfectly disentangled,

the empirical variance in the dimension corresponding to the fixed factor will be 0. The representations are normalised so that the arg min is invariant to rescaling of the representations in each dimension. Since both inputs and outputs lie in a discrete space, the optimal classifier is the majority-vote classifier, and the metric is the error rate of the classifier. Lower values imply better disentanglement. 788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

- 3. Mutual Information Gap (MIG) (Chen et al., 2018): The difference between the top two latent variables with the highest mutual information. Empirical mutual information between a latent representation z_i and a ground truth factor v_k , is estimated using the joint distribution defined by $q(z_j, v_k) =$ $\sum_{n=1}^{N} p(v_k) p(n|v_k) q(z_j|n).$ A higher mutual information implies that z_i contains a lot of information about v_k , and the mutual information is maximal if there exists a deterministic, invertible relationship between z_i and v_k . MIG values are in the interval [0, 1], with higher values implying better disentanglement.
- 4. *Modularity* (Ridgeway and Mozer, 2018): The deviation from an ideally modular case of latent representation. If latent vector dimension *i* is ideally modular, it will have high mutual information with a single factor and zero mutual information with all other factors. A deviation δ_i of 0 indicates perfect modularity and 1 indicates that this dimension has equal mutual information with every factor. Thus, $1 - \delta_i$ is used as a modularity score for vector dimension i and the mean of $1 - \delta_i$ over *i* as the modularity score for the overall representation. Higher values imply better disentanglement.
- 5. *Explicitness* (Ridgeway and Mozer, 2018): Mean of the ROC area-under-the-curve (AUC_{jk}) of a one-versus-rest logistic-regression classifier that takes the latent vectors as input and has factor values as targets, over a factor index j and an index k on values of factor j. Represents the coverage of the representation, in other words, how well each factor is represented. Higher values imply better disentanglement.
- 6. Disentanglement Score (Eastwood and

Williams, 2018): The degree to which a representation factorises or disentangles the underlying factors of variation, with each variable (or dimension) capturing at most one generative factor. It is computed as a weighted average of a disentanglement score $D_i = (1 - H_K(P_i))$ for each latent dimension variable c_i , on the relevance of each c_i , where $H_K(P_i)$ denotes the entropy and P_{ij} denotes the 'probability' of c_i being important for predicting z_i . If c_i is important for predicting a single generative factor, the score will be 1. If c_i is equally important for predicting all generative factors, the score will be 0. Higher values imply better disentanglement.

837

838

841

842

843

845

852

855

856

864

867

870

871

872

873

- 7. Completeness Score (Eastwood and Williams, 2018): The degree to which each underlying factor is captured by a single latent dimension variable. For a given z_j it is given by $C_j = (1 H_D(\tilde{P}.j))$, where $H_D(\tilde{P}.j) = -\sum_{d=0}^{D-1} \tilde{P}_{dj} log_D \tilde{P}_{ij}$ denotes the entropy of the $\tilde{P}.j$ distribution. If a single latent dimension variable contributes to z_j 's prediction, the score will be 1 (complete). If all code variables contribute equally to z_j 's prediction, the score will be 0 (maximally over-complete). Higher values imply better disentanglement.
- 8. Informativeness Score (Eastwood and Williams, 2018): The amount of information that a representation captures about the underlying factors of variation. Given a latent representation c, It is quantified for each generative factor z_j by the prediction error $E(z_j, \hat{z}_j)$ (averaged over the dataset), where E is an appropriate error function and $\hat{z}_j = f_j(c)$. Lower values imply better disentanglement.

C Further Experimental Results

UMAP plot Alternative dimensionality reduction method, used to visualise the clustering of DSR patterns, as seen in Figure 6.



Figure 6: UMAP plot of 9370 definition representations (128 dimensions), generated from Unsupervised VAE (U), DSR supervision (S) and Conditional VAE (C).

D Source code

The complete experimental pipeline is available as supplementary software for this paper (code.7z), and should be soon available to the public. 875 876

877

878

879

880

881