## **DIFFUCOMET: Contextual Commonsense Knowledge Diffusion**

Anonymous ACL submission

#### Abstract

001 Inferring contextually-relevant and diverse commonsense to understand narratives remains challenging for knowledge models. In this work, we develop a series of knowledge models, DIFFUCOMET, that leverage diffusion to learn to reconstruct the implicit semantic connections between narrative contexts and relevant commonsense knowledge. Across multiple diffusion steps, our method progressively refines a representation of commonsense facts that is anchored to a narrative, producing contextually-relevant and diverse commonsense inferences for an input context. To evaluate DIFFUCOMET, we introduce new metrics for commonsense inference that more closely measure knowledge diversity and contextual relevance. Our results on two different 017 benchmarks, ComFact and WebNLG+, show that knowledge generated by DIFFUCOMET 020 achieves a better trade-off between commonsense diversity, contextual relevance and align-021 ment to known gold references, compared to 022 baseline knowledge models.

## 1 Introduction

026

027

Identifying the commonsense inferences that underlie narratives, such as stories or dialogues (Guan et al., 2019; Zhou et al., 2022), is crucial to understanding those same narratives. For example, to understand why "Hank ... got the shopping bags" in the context in Figure 1, a model would need to infer that (1) Hank was not finished wrapping gifts, and so (2) would need to buy more wrapping paper. However, comprehensively inferring these diverse, yet implicit, commonsense inferences that are relevant to a context remains a challenging task.

Recent methods for identifying contextuallyrelevant commonsense inferences (Bosselut et al., 2021; Tu et al., 2022; Peng et al., 2022) use knowledge models (Bosselut et al., 2019; West et al., 2022) to generate commonsense facts. While



Figure 1: Overview of our diffusion-based contextual commonsense knowledge generation.

knowledge models have been less brittle than previous retrieval-based methods for commonsense inference, they have two major shortcomings. First, they are trained to verbalize tuples from general commonsense knowledge graphs (Sap et al., 2019; Hwang et al., 2021), leading them to produce valid, but often contextually-irrelevant, commonsense inferences when applied out-of-the-box to real narratives. Second, because they are trained using autoregressive training objectives, they subsequently decode high-likelihood, non-diverse sequences that only identify limited collections of commonsense inferences relevant to an input context.

In this work, we address these challenges of contextual commonsense knowledge generation by developing **Diffu**sion (Ho et al., 2020) **COM**monsEnse Transformer (Bosselut et al., 2019) models. DIFFUCOMET models (shown in Figure 1) uses diffusion-based decoding to generate relevant knowledge embeddings that are constrained to the narrative context. Over multiple it-

erations of constrained diffusion, our models refine a latent representation of the semantic connections between a context and its relevant facts, ensuring that it generates commonsense knowledge that is more contextually relevant to the narrative. At the same time, by jointly refining multiple fact embeddings during diffusion, DIFFUCOMET also generates more diverse inferences than comparable-size autoregressive knowledge models.

062

063

064

067

076

080

086

089

095

097

100

101

102

103

104

105

106

108

109

110

We evaluate DIFFUCOMET models using traditional NLG metrics (e.g., BLEU; Papineni et al., 2002) commonly used for evaluating knowledge models. However, these metrics focus on surface form matching to gold references, and fall short of measuring the diversity of commonsense inferences and their semantic relevance to real narrative contexts. Our second contribution is a novel set of metrics that assess the diversity and contextual relevance of knowledge generated by knowledge models. Using both the traditional evaluation metrics and our new suite, we evaluate our models on a commonsense inference linking benchmark (Gao et al., 2022a) that covers both social and physical knowledge, and a second knowledge generation benchmark that involves extracting RDF triplets from language, WebNLG+ (Ferreira et al., 2020).

Our result show that DIFFUCOMET models generate knowledge that achieves a better balance of diversity and contextual relevance compared to other knowledge models. DIFFUCOMET models also more robustly generalize to generate knowledge for out-of-distribution narratives, and are better at producing novel knowledge tuples that are not in their initial training set. Finally, on our second benchmark, WebNLG+, we verify that our diffusion modeling method also generalizes well to a completely new factual knowledge generation task beyond the commonsense domain.

## 2 Background: Diffusion Models

Diffusion models learn to construct synthetic data from random noise. They use a forward process to gradually corrupt real data samples with additive noise, and learn a reverse process to recover (or de-noise) the corrupted data samples. Through the de-noising of corrupted data, diffusion models learn to map from a random noise distribution to their target data distribution, which grounds their synthetic data generation.

In this paper, we adopt the  $DDPM^1$  (Ho et al.,

2020) formulation of the forward and reverse diffusion processes. Specifically, based on a sample  $z_0$ from a continuous input data distribution  $q(z_0)$ , the forward process constructs noisy sample  $z_t$  over a sequence of time steps  $t \in \{1, 2, ..., T\}$ . In DDPM,  $z_t$  is sampled from a Gaussian distribution conditioned on the previous sample  $z_{t-1}$ , given by:

$$q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I}) \quad (1)$$

111

112

113

114

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

where  $\beta_t$  is a noise schedule hyperparameter unique to each diffusion step.

In the reverse process, diffusion models learn an inverse distribution  $q(\mathbf{z}_{t-1}|\mathbf{z}_t)$  to de-noise samples created by the forward process. To more precisely couple the intermediate states of the reverse process with the final de-noised sample  $\mathbf{z}_0$ , Diffusion-LM (Li et al., 2022) reformulates the task of predicting  $\mathbf{z}_{t-1}$  as directly predicting  $\mathbf{z}_0$  (based on  $\mathbf{z}_t$ ), and uses a mean-squared error training loss on the  $\mathbf{z}_0$ prediction at each time step<sup>2</sup>:

$$\mathcal{L}_{z_0\text{-}mse} = \sum_{t=1}^{T} \mathbb{E} \|\mathbf{z}_0 - f_\theta(\mathbf{z}_t, t)\|^2 \qquad (2)$$

where  $f_{\theta}(\mathbf{z}_t, t) = \hat{\mathbf{z}}_0^{t-1}$  denotes the model's learned prediction of  $\mathbf{z}_0$  at the reverse stage of step t to t-1. To formulate  $\hat{\mathbf{z}}_0^{t-1}$  as a refinement of the former reverse stage's output  $\hat{\mathbf{z}}_0^t$ , Bit-Diffusion (Chen et al., 2022) improves the model function of predicting  $\mathbf{z}_0$  with self-conditioning, *i.e.*,  $\hat{\mathbf{z}}_0^{t-1} = f_{\theta}(\hat{\mathbf{z}}_0^t, \mathbf{z}_t, t)$ . At inference time, the noisy sample at step t is predicted from  $\hat{\mathbf{z}}_0^t$  via the Eq. (1) forward process, denoted as  $\hat{\mathbf{z}}_t$  to replace the unknown gold input  $\mathbf{z}_t$ , while the initial input  $\mathbf{z}_T$  is pure Gaussian noise sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .

## **3** Contextual Knowledge Diffusion

In this section, we first introduce the task of contextual commonsense knowledge generation, and then propose DIFFUCOMET, our diffusion approach for this task. The overview of our method is presented in Figure 1.

**Task Description** Given a narrative sample S as context, *e.g.*, a story snippet or a dialogue, the model needs to generate commonsense inferences as a set facts  $\mathcal{K} = \{k_1, ..., k_n, ..., k_N\}$ , which are relevant for understanding the situation described in the context. Each fact  $k_n = (h_n, r_n, a_n)$  is represented as a triple containing a head entity  $h_n$ , a tail

<sup>&</sup>lt;sup>1</sup>Denoising Diffusion Probabilistic Models

<sup>&</sup>lt;sup>2</sup>We include more detailed formulation of the reverse diffusion training in Appendix A.



Figure 2: Knowledge diffusion based on facts or entities. Dashed arrows denote the forward process used for constructing gold references at the training phase. Solid arrows denote the reverse process used for generating knowledge with attention to the narrative context.

(attribute) entity  $a_n$ , and a relation  $r_n$  connecting them, *e.g.*, (*wrapping paper*, *used for*, *wrap gifts*), as shown in Figure 1. We denote the set of unique head entities, relations and tail entities in  $\mathcal{K}$  as  $\mathcal{H}$ ,  $\mathcal{R}$  and  $\mathcal{A}$ , respectively.

155

156

157

160

161

162

163

164

165

168

169

170

171

172

174

**Contextualization** We ground knowledge diffusion on the given context S by using encoderdecoder cross attention, inspired by SeqDiffuSeq (Yuan et al., 2022). In particular, we use a BART (Lewis et al., 2020) encoder  $f_{\theta_s}$  to learn the context encoding that represents S as hidden state  $z_S$ :

$$\mathbf{z}_{\mathcal{S}} = f_{\theta_s}(\mathcal{S}) \tag{3}$$

Then, a BART decoder  $f_{\theta_z}$ , serving as the diffusion module, learns to predict the de-noised data sample  $z_0$ . Given the context hidden state  $z_S$  (via cross-attention to the encoder  $f_{\theta_s}$ ),  $f_{\theta_z}$  makes a prediction of  $z_0$  at time step t-1 (*i.e.*,  $\hat{z}_0^{t-1}$ ) based on its former prediction  $\hat{z}_0^t$  and time step t's noisy sample  $z_t$ :

$$\hat{\mathbf{z}}_0^{t-1} = f_{\theta_z}(\hat{\mathbf{z}}_0^t, \mathbf{z}_t, t | \mathbf{z}_{\mathcal{S}})$$
(4)

175 Unlike traditional transformer decoders (Vaswani 176 et al., 2017), the diffusion module  $f_{\theta_z}$  applies a 177 bi-directional self-attention to  $\hat{\mathbf{z}}_0^t$  and  $\mathbf{z}_t$ , since all positions of  $\hat{\mathbf{z}}_0^{t-1}$  are decoded simultaneously, *i.e.*, in non-autoregressive manner.<sup>3</sup>

178

179

180

181

182

184

185

186

188

189

192

193

194

195

196

198

199

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

**Discrete Knowledge Diffusion** We consider two formulations for representing discrete knowledge in continuous embedding spaces for diffusion: **DIF-FUCOMET-Fact**, where we learn to reconstruct continuous representations of facts  $k_n$  using diffusion, and **DIFFUCOMET-Entity**, where we use separate diffusion processes to reconstruct head  $h_n$  and tail  $a_n$  representations and then predict the relation between them to complete the fact. We highlight these differences in Figure 2.

For diffusion on the fact-level embedding space (**DIFFUCOMET-Fact**), we first pre-train a BART encoder  $f_{\theta_e}$  to produce an embedding  $\mathbf{e}_n$  of each fact  $k_n$  in the knowledge set  $\mathcal{K}$  (with embedding size d same as the hidden state size of BART):

$$\mathbf{e}_n = f_{\theta_e}(k_n) \in \mathbb{R}^d \tag{5}$$

where we input the concatenation of each fact's head, relation and tail tokens to the encoder  $f_{\theta_e}$ , and take the output hidden state of a start token  $\langle s \rangle$  as the embedding of the fact. The initial input  $\mathbf{z}_0$  of the forward diffusion process is then sampled from a Gaussian centered on the concatenation of all fact embeddings  $\mathbf{e} = [\mathbf{e}_1; \mathbf{e}_2; ...; \mathbf{e}_{|\mathcal{K}|}] \in \mathbb{R}^{d \times |\mathcal{K}|}$ , formulated as  $q_e(\mathbf{z}_0 | \mathbf{e}) = \mathcal{N}(\mathbf{z}_0; \mathbf{e}, \beta_0 \mathbf{I})$ .

In the reverse process, the diffusion module  $f_{\theta_z}$ is trained to generate the final output  $\hat{\mathbf{z}}_0^0$  (using time step 1's input  $\mathbf{z}_1$  and  $\hat{\mathbf{z}}_0^1$ ) as its predicted fact embeddings  $\hat{\mathbf{e}}$ , *i.e.*,  $\hat{\mathbf{e}} = \hat{\mathbf{z}}_0^0 = f_{\theta_z}(\hat{\mathbf{z}}_0^1, \mathbf{z}_1, 1 | \mathbf{z}_S)$ . Finally, we pre-train another BART decoder  $f_{\theta_g}$  to generate the synthetic fact  $\hat{k}_n$  with conditioned on the diffusion module's predicted *n*-th embedding  $\hat{\mathbf{e}}_n = \hat{\mathbf{e}}[:][n], (n = 1, 2, ..., |\mathcal{K}|)^4$ :

$$\hat{k}_n = f_{\theta_g}(\hat{\mathbf{e}}_n) \tag{6}$$

For diffusion on the entity-level embedding space (**DIFFUCOMET-Entity**), we use a pipeline to generate head entities, tail entities and their relations. First, to generate head entities, we use a similar process as in **DIFFUCOMET-Fact**, *i.e.*, pretrain a BART encoder to produce a gold embedding of each unique head entity  $h_i \in \mathcal{H}$  (for training the

<sup>&</sup>lt;sup>3</sup>More implementation details of the diffusion module  $f_{\theta_z}$  are presented in Appendix B.1.

<sup>&</sup>lt;sup>4</sup>At inference time, the maximum value of n (number of generated facts) can be arbitrary depending on the user's choice. In Appendix B.2, we introduce how we control the number of facts that our models generate for each context.

220diffusion module), and then pre-train a BART de-221coder to generate synthetic head entities  $\hat{h}_i$  from222the diffusion module's predicted embeddings. Each223predicted head entity  $\hat{h}_i$  is then appended to the224context (*i.e.*, S in Eq. 3), expanding the context to225 $S_i = [S, \hat{h}_i]$ . A second diffusion module predicts226embeddings of synthetic tail entities  $\hat{a}_j$  related to227 $S_i$  (trained using gold embeddings of tail entities228 $a_j \in A$  that possess relations  $r_{ij} \in \mathcal{R}$  to the gold229head  $h_i$ ). A final BART model predicts the relation230 $\hat{r}_{ij}$  between each pair of generated head and tail231entities, grounded on the context.

**Embedding Module Training** We pretrain the embedding modules  $(f_{\theta_e}, f_{\theta_g})$ , which focus on modeling generic knowledge representations independent to the context, before the diffusion modules  $(f_{\theta_s}, f_{\theta_z})$ , which learn the specific mapping from the context to its relevant knowledge. When training the diffusion modules, we freeze the pretrained embedding modules.

240

241

242

243

244

245

247

248

250

251

257

258

260

261

264

To pretrain the fact (or entity) embedding modules, we minimize the decoder's negative loglikelihood of re-constructing facts k (or entity h or a) in the full set of knowledge  $\mathcal{K}_{full}$  involved in the whole narrative dataset (or domain), based on its embedding given by encoder  $f_{\theta_e}$ :

$$\mathcal{L}_{\theta_e,\theta_g} = -\log p_{\theta_g}(k|f_{\theta_e}(k)) \tag{7}$$

**Diffusion Module Training** We optimize a dual loss to train the diffusion modules. First, we consider the mean-square error loss of the diffusion module's de-noised sample prediction  $\hat{\mathbf{z}}_0^t$  at each time step t, compared to the reference sample  $\mathbf{z}_0$ (for t > 0) and gold embeddings e (for t = 0):

$$\mathcal{L}_{\theta_s,\theta_z}^{mse} = \mathbb{E} \| \mathbf{e} - \hat{\mathbf{z}}_0^0 \|^2 + \sum_{t=1}^{T-1} \mathbb{E} \| \mathbf{z}_0 - \hat{\mathbf{z}}_0^t \|^2 \quad (8)$$

We also use an anchor loss (Gao et al., 2022b) to supervise the final fact (or entity) generation. For each time step t, we minimize the negative loglikelihood of the embedding module decoder (with frozen parameters  $\theta_g$ ) generating each fact  $k_n$  in knowledge set  $\mathcal{K}$ , based on the diffusion module's predicted de-noised sample  $\hat{\mathbf{z}}_0^t$ :

$$\mathcal{L}_{\theta_s,\theta_z}^{gen} = \sum_{t=0}^{T-1} \sum_{n=1}^{|\mathcal{K}|} -\log p_{\theta_g}(k_n | \hat{\mathbf{z}}_0^t[:][n]) \quad (9)$$

where  $\hat{\mathbf{z}}_0^t[:][n]$  is the predicted de-noised representation of  $k_n$ . The final loss is  $\mathcal{L}_{\theta_s,\theta_z} = \mathcal{L}_{\theta_s,\theta_z}^{mse} + \gamma \mathcal{L}_{\theta_s,\theta_z}^{gen}$ , where  $\gamma$  is a tunable hyperparameter. **Inference** At inference time, the reverse diffusion process is initialized with noise sampled from the Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , while the embedding module encoder  $f_{\theta_e}$ , which provides gold diffusion references for training, is not used.

### 4 Evaluation

Prior work in commonsense knowledge generation (Hwang et al., 2021; Da et al., 2021) evaluated knowledge models using traditional NLG metrics (e.g., BLEU; Papineni et al., 2002) in controlled studies with KGs, where the inputs to the models were head entities and relations and the knowledge model produced tail attributes. In practice, however, knowledge models are used to generate implicit commonsense inferences for natural language contexts (Ismayilzada and Bosselut, 2023), requiring generated inferences to be relevant to a more complex input than a basic KG head entity, and necessitating diverse generated inferences that comprehensively augment the context. However, traditional NLG metrics fall short of measuring these important dimensions because they measure surface form overlap between model outputs and references, which rewards generating facts with similar or duplicated semantics, limiting diversity.

Motivated by these shortcomings, we propose novel evaluation metrics that assess the diversity and contextual relevance of generated knowledge. First, to eliminate the effect of knowledge repetition in generations, we cluster similar facts and treat each fact cluster (instead of each single fact) as a unit piece of knowledge. In particular, we use the DBSCAN (Ester et al., 1996) algorithm to group gold facts  $\mathcal{K} = \{k_1, k_2, ..., k_N\}$  and generated facts  $\hat{\mathcal{K}} = \{\hat{k}_1, \hat{k}_2, ..., \hat{k}_{\hat{N}}\}$  into clusters  $\mathcal{C} = \{c_1, c_2, ..., c_M\}$  and  $\hat{\mathcal{C}} = \{\hat{c}_1, \hat{c}_2, ..., \hat{c}_{\hat{M}}\}$ , respectively. We test two methods for measuring the similarity of facts for clustering: word-level edit distance (Levenshtein et al., 1966), which measures the difference of two facts' surface-form tokens, and Euclidean distance of Sentence-BERT (Reimers and Gurevych, 2019) embeddings, which measures the semantic difference of two facts. Based on these clusters, we develop three metrics to measure the diversity of generated facts, their contextual relevance, and their alignment to gold references, as shown in Figure 3.

**Diversity.** To measure the diversity of generated facts (*i.e.*, amount of distinctive knowledge being generated), we count the number of fact clusters

265

266

267

271 272

270

274 275 276

277

278

279

280

281

283

284

285

287

290

291

292

293

294

296

297

299

300

301

302

304

305

306

308

309

310

311

312

313

314



Figure 3: Illustration of clustering-based evaluation metrics for contextual commonsense knowledge generation.

315(# Clusters), *i.e.*,  $\hat{M}$  (or M for gold references).316We also report the number of facts (# Facts), *i.e.*,  $\hat{N}$ 317(or N for gold references), to compare the number318of fact clusters to the number of generated facts319produced by the models.

320

321

326

328

332

335

336

337

341

342

344

345

**Relevance.** We measure the relevance of the fact clusters to the narrative context using a fact linker<sup>5</sup> trained on the *ComFact* dataset (Gao et al., 2022a) that scores the relevance of each fact  $\hat{k}_n$  to the context S, denoted as  $rel(\hat{k}_n, S) \in [0, 1]$ . The relevance score of a fact cluster  $\hat{c}_m$  is defined as the average relevance score of its facts, *i.e.*,  $\sum_{\hat{k}_n \in \hat{c}_m} rel(\hat{k}_n, S)/|\hat{c}_m|$ . Finally, we measure the average relevance over all fact clusters in  $\hat{C}$ :

$$rel(\hat{\mathcal{C}}, \mathcal{S}) = \frac{1}{\hat{M}} \sum_{\hat{c}_m \in \hat{\mathcal{C}}} \frac{1}{|\hat{c}_m|} \sum_{\hat{k}_n \in \hat{c}_m} rel(\hat{k}_n, \mathcal{S})$$
(10)

We note that **Relevance** can be viewed as a *precision* measure for generated facts, which tends to decrease as more facts are generated because irrelevant facts are more likely to be generated.

Alignment measures the average similarity of generated facts to gold fact clusters. Specifically, we define a function  $sim(\hat{k}_i, k_j) \in [0, 1]$  to measure the pairwise similarity between a generated fact and a gold reference (using similar distance functions to define clusters above<sup>6</sup>). Using this function, we measure the maximum pairwise similarity of generated facts to references in each gold cluster  $c_m \in C$ , which serves as the alignment score to the gold cluster. Finally, we average the alignment scores of generated facts to all gold clusters:

$$sim(\hat{\mathcal{K}}, \mathcal{C}) = \frac{1}{M} \sum_{c_m \in \mathcal{C}} \max_{\substack{\hat{k}_i \in \hat{\mathcal{K}}, \\ k_j \in c_m}} sim(\hat{k}_i, k_j) \quad (11)$$

We note that **Alignment** can be viewed as the generated facts' *recall* of gold fact clusters, which tends to increase as more facts are generated because more facts will be aligned to gold clusters. Given this trade-off between Relevance and Alignment, we also present the harmonic mean of Relevance and Alignment as an overall evaluation of the two dimensions, denoted as **RA-F1**. 346

347

348

350

351

352

354

355

356

357

358

359

361

362

363

364

365

367

369

370

371

372

373

375

376

377

379

381

382

## **5** Experimental Settings

**Datasets** First, we evaluate our approach on the ComFact (Gao et al., 2022a) benchmark, where models need to generate ATOMIC $_{20}^{20}$  (Hwang et al., 2021) social commonsense facts that are relevant to narrative contexts sampled from four diverse corpora: PERSONA-CHAT (Zhang et al., 2018), Mu-Tual (Cui et al., 2020), ROCStories (Mostafazadeh et al., 2016) and CMU Movie Summary (Bamman et al., 2013). We only use training data from the ROCStories portion of ComFact, to enable the evaluation of zero-shot generalization on the other three partitions of the dataset. Our fact embedding module is pretrained on the full ATOMIC $^{20}_{20}$ knowledge base, which contains  $\sim 972K$  commonsense facts after preprocessing.<sup>7</sup> We also evaluate our approach in a conceptually different setting, the WebNLG+ 2020 (Ferreira et al., 2020) dataset, which consists of RDF (Ora, 1999) facts sampled from the DBpedia (Lehmann et al., 2015) knowledge base, with corresponding natural language texts verbalizations. The task is to generate the sampled RDF facts given their verbalizations. We use  $\sim 13$ k facts from the training data to pretrain the fact embedding module.

**Baselines** We train DIFFUCOMET using BARTbase and BART-large as pretrained models, and compare with three baselines developed on the same backbones: a) a **Greedy** baseline that is

<sup>&</sup>lt;sup>5</sup>Fact linking models predict the relevance of knowledge tuples to textual passages (Gao et al., 2022a)

<sup>&</sup>lt;sup>6</sup>Further details on exact definitions are in Appendix C.1.

<sup>&</sup>lt;sup>7</sup>More data preprocessing details are in Appendix D.

Model	# Facts	# Clusters	Relevance	Alignment	RA-F1	BLEU	METEOR	<b>ROUGE-L</b>
Greedy-COMET	1.96	1.19	61.42	50.64	55.51	18.01	52.32	54.96
Sampling-COMET	15.00	8.39	56.19	77.97	65.31	12.69	44.43	45.58
Beam-BART	15.00	4.60	64.35	71.35	67.67	13.11	47.70	46.35
Beam-COMET	15.00	5.09	65.03	71.64	68.18	16.97	47.39	47.19
Grapher	5.08	2.60	68.29	40.58	50.91	1.40	23.96	27.21
DIFFUCOMET-Fact	12.88	5.24	65.64	71.65	<u>68.51</u>	15.98	<u>50.06</u>	<u>51.44</u>
DIFFUCOMET-Entity	12.89	<u>5.67</u>	<u>66.39</u>	<u>74.38</u>	70.16	<u>17.01</u>	47.61	48.40
Gold	10.55	5.64	80.90	-	-	-	-	-

Table 1: Evaluation results on the ROCStories portion of ComFact. Both DIFFUCOMET models presented are developed based on BART-large. Models with suffix "-COMET" and "-BART" are fine-tuned on COMET-BART and BART-large. Presented results of our proposed metrics are based on fact clustering *w.r.t.* embedding Euclidean distance. Best and second-best results (excluding Gold references) are **bolded** and <u>underlined</u>, respectively.

trained to autoregressively generate the concatenation of all relevant facts,<sup>8</sup> b) a **Sampling** baseline that uses nucleus sampling (Holtzman et al., 2019) to generate multiple individual facts in parallel, and c) a Diverse Beam search baseline that uses diverse beam search to generate multiple inferences in parallel. We also compare our models trained using BART-large to baselines developed on models of similar scale: d) the aforementioned greedy decoding, sampling and beam search baselines trained from COMET-BART (Hwang et al., 2021), a BART-large model further pre-trained on  $\text{ATOMIC}_{20}^{20}$  for commonsense knowledge completion, and e) Grapher (Melnyk et al., 2022), which trains a T5-large (Raffel et al., 2020) model to generate entities (nodes) related to the context, followed by a MLP classifier to predict the relations (edges) between entities.

386

393

396

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

Metrics We evaluate these methods on our clustering-based metrics described in Section 4. As the clustering algorithm (i.e., DBSCAN) used in our metrics has an adjustable clustering granularity controlled by a distance threshold, we consider a range of distance thresholds and take the average of evaluation results across all thresholds in the range, allowing us to avoid biasing our metrics to a specific distance threshold.<sup>9</sup> For ComFact, we also test on the metrics from Hwang et al., 2021 for evaluating commonsense knowledge generation, including BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005) and ROUGE-L (Lin, 2004). For evaluation on WebNLG+ 2020, we also report the official metrics for this dataset's challenge (Ferreira et al., 2020), which construct

optimal pairings between predicted facts and gold references, and then compute precision, recall, and F1 scores based on the surface-form matching of paired facts. We denote these WebNLG metrics as **Web-Prec.**, Web-Rec. and Web-F1.<sup>10</sup> 417

418

419

420

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

#### 6 Results and Analysis

Table 1 shows evaluation results on the ROCStories portion of the ComFact benchmark for our DIFFUCOMET models developed based on BARTlarge.<sup>11</sup> On our new cluster-based metrics, DIFFU-COMET models demonstrate a better balance between diversity and accuracy in contextual knowledge generation. Specifically, DIFFUCOMET models achieve Relevance and Alignment scores that are both comparable to the best baseline results, contributing to their higher overall RA-F1 measures, while also producing a larger number of distinct knowledge clusters. By contrast, the Greedy, Sampling and Grapher baselines significantly sacrifice one or two dimensions of diversity and quality w.r.t. # Clusters, Relevance and Alignment. Beam baselines consistently underperform DIFFUCOMET on cluster metrics.

For evaluation on the traditional NLG metrics, we find that DIFFUCOMET models score higher overall than most baseline models on metrics that check the alignment with gold references, *i.e.*, BLEU, METEOR and ROUGE-L, except for the Greedy decoding baseline, whose higher scores are artificially high because it generates very little knowledge, *i.e.*, only  $\sim$ 2 facts per context. We also include further comparisons of models with

<sup>&</sup>lt;sup>8</sup>Facts are concatenated by a special token *<fsep>*.

<sup>&</sup>lt;sup>9</sup>We include more details of our clustering threshold selection in Appendix C.2.

<sup>&</sup>lt;sup>10</sup>More details of WebNLG metrics are in Appendix C.3.

<sup>&</sup>lt;sup>11</sup>Presented results of our metrics are based on fact clustering *w.r.t.* embedding Euclidean distance. Results based on word-level edit distance are included in Appendix F.1, and promote the same conclusions.

Model	Validity	Relevance
Sampling-COMET	49.45	30.20
Beam-COMET	<b>74.80</b>	42.81
DIFFUCOMET-Fact	70.00	<u>48.27</u>
DIFFUCOMET-Entity	<u>74.15</u>	<b>54.18</b>
Gold	94.79	82.04

Table 2: Human evaluation results.

Model	# Novel Facts	# Novel Clusters
Sampling-COMET	0.26	0.19
Beam-COMET	0.27	0.17
DIFFUCOMET-Fact	0.30	<u>0.20</u>
DIFFUCOMET-Entity	0.30	<b>0.24</b>

Table 3: Novelty of generated knowledge.

BART-*base* backbones in Appendix F.1, where our models outperform baselines by a larger gap, *i.e.*,  $\sim 15\%$  absolute RA-F1 improvement on average.

We also test DIFFUCOMET's ability to generalize to out-of-domain contexts using the other portions of ComFact with contexts sampled from PersonaChat, MuTual and MovieSummaries. We report generalization results to the above three portions in Appendix Tables 8-13, and observe similar results where DIFFUCOMET-Entity outperforms baselines by ~5% RA-F1 and produces ~20% more knowledge clusters.

The results of our automatic evaluation are also supported by our human evaluation. We hire Amazon Mechanical Turk workers<sup>12</sup> to evaluate the validity and contextual relevance of models' generated knowledge on the ROCStories portion of ComFact. Specifically, given a narrative context and a list of commonsense facts that a model generates about the context, we ask three workers to independently judge whether each fact is valid and relevant<sup>13</sup> to the context, and take their majority vote as the assessment. In Table 2, we see that DIF-FUCOMET models produce *valid* facts at about the same rate as the best baseline, but produce facts that are far more relevant to the narrative context.

475 Novelty DIFFUCOMET models also produce
476 more novel commonsense inferences. A histori477 cal advantage of knowledge models (*e.g.*, COMET)
478 was their ability to generate knowledge beyond the
479 graphs they used for pretraining (Bosselut et al.,
480 2019), making them powerful tools to generate



Figure 4: DIFFUCOMET performance at different diffusion steps during inference. Both DIFFUCOMET-Fact and DIFFUCOMET-Entity are developed based on BART-large and tested on the ROCStories portion of ComFact. Beam-COMET performance is shown as a baseline, with the number of decoded facts set to match DIFFUCOMET-Entity at each diffusion step.

commonsense knowledge for unseen narratives. To test the novelty of generated commonsense knowledge from DIFFUCOMET, we develop a heuristic method that identifies knowledge as *novel* if its maximum pair-wise (Sentence-BERT embedding) cosine similarity to *ComFact* gold references is lower than 0.45. However, as this cut-off would likely cause invalid and irrelevant facts to be considered novel, we only include facts whose relevance score is higher than 0.97.<sup>14</sup> In Table 3, we see that DIFFUCOMET models produce more novel facts and clusters compared to baselines.<sup>15</sup>

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

**Diffusion Steps** To investigate how DIFFU-COMET's multiple rounds of knowledge representation refinement through the diffusion process affect the quality of generated knowledge, we record the performance of our DIFFUCOMET models as a function of diffusion steps conducted during inference. Figure 4 shows how DIFFUCOMET's performance varies when knowledge is generated at earlier time steps.

<sup>&</sup>lt;sup>12</sup>Details on workers and their payment are in Appendix E

<sup>&</sup>lt;sup>13</sup>*invalid* facts are automatically labeled *irrelevant* 

<sup>&</sup>lt;sup>14</sup>Thresholds are tuned by a manual check of 100 sampled results to ensure a decent cutoff of novel and relevant facts.

<sup>&</sup>lt;sup>15</sup>We conduct analysis on some examples of novel facts in a case study in Appendix F.3.

Model	Web-Prec.	Web-Rec.	Web-F1
Beam-BART	73.36	76.27	74.75
Grapher	71.20	73.00	71.90
DIFFUCOMET-Fact	<u>76.30</u>	<u>78.07</u>	<u>77.19</u>
DIFFUCOMET-Entity	<b>80.68</b>	82.89	<b>81.74</b>

Table 4: Results on **WebNLG+ 2020**. Official metrics used for the benchmark challenge are presented.

522

523

529

530

531

535

537

## We find that DIFFUCOMET models gradually produce more facts and more diverse facts (i.e., # Clusters) as the number of diffusion steps increase, indicating that the multiple rounds of diffusion produce a more separable representation capable of representing more facts. While the greater number of facts leads to a slight drop in contextual relevance across the generated facts, a greater corresponding increase in alignment to the gold clusters (as observed by the increase in Alignment and RA-F1) also emerges. On RA-F1, DIFFUCOMET-Fact surpasses Beam-COMET<sup>16</sup> as the diffusion steps increase to larger than 200, and DIFFUCOMET-Entity consistently scores higher and continues benefiting from further diffusion, even after 1000 diffusion steps. These results shows that multi-step refinement of facts via diffusion effectively improves contextual knowledge generation.

## 6.1 WebNLG+ 2020 Benchmark

Finally, to test whether our method generalizes outside the domain of generating commonsense inferences, we present our evaluation results on the WebNLG+ 2020 dataset in Table 4. DIFFU-COMET models achieve better performances on the WebNLG factual knowledge generation task, verified by the official metrics of the benchmark.<sup>17</sup> This results suggests that our diffusion approach to knowledge graph construction could be adapted to other knowledge generation tasks.

## 7 Related Work

**Commonsense Knowledge Grounding** To augment NLP systems with commonsense knowledge, various systems for question answering (Zhang et al., 2022; Yasunaga et al., 2021, 2022) and narrative generation (Ji et al., 2020; Zhou et al., 2022) use retrieval methods based on heuristics to link rel-

evant facts from commonsense knowledge graphs (Speer et al., 2017; Sap et al., 2019; Gao et al., 2023). However, these systems typically have low precision when adapted to more general and complex commonsense linking (Hwang et al., 2021; Jiang et al., 2021). Gao et al., 2022a developed commonsense fact linking to improve retrieval precision, but this requires inefficiently traversing all candidate facts to check their contextual relevance. 538

539

540

541

542

543

544

545

546

547

548

549

551

552

553

554

555

556

557

558

559

561

562

564

565

566

567

569

570

571

572

573

574

575

576

578

579

580

581

582

583

584

586

Due to above limitations of retrieval-based knowledge grounding, one line of research (Bosselut et al., 2021; Tu et al., 2022) uses knowledge models (Bosselut et al., 2019; West et al., 2022) to generate tail inferences from narrative statements. However, these methods often produce irrelevant facts as the knowledge models are pre-trained for context-free knowledge graph completion. Finally, developing new knowledge models to learn contextual commonsense generation turns out to be a promising track of research, while current works are limited to simple physical (Zhou et al., 2022) or RDF-style factual (Melnyk et al., 2022) knowledge. We build new models to address contextual commonsense generation in a more general scope.

**Diffusion Models** Considerable recent works (Gao et al., 2022b; Lin et al., 2022; Han et al., 2024) have developed methods to improve text generation with diffusion models (Sohl-Dickstein et al., 2015; Song and Ermon, 2019; Ho et al., 2020). However, the potential of diffusion models in text-to-knowledge generation is still under-explored. In this paper, we introduce diffusion models for the task of contextual knowledge generation.

## 8 Conclusion

In this work, we leverage the power of diffusion models for contextual commonsense knowledge generation, and formulate novel metrics to highlight important dimensions of diversity and contextual relevance for this task. Our diffusion knowledge models, DIFFUCOMET, outperform various autoregressive knowledge models, producing more diverse, novel, and contextually-relevant commonsense knowledge, and achieving better out-ofdistribution performance. Finally, our analysis reveals how DIFFUCOMET refines implicit knowledge representations over the course of the diffusion process to produce more relevant and diverse inferences, hinting at our method's potential benefit in other text-to-graph generation tasks.

<sup>&</sup>lt;sup>16</sup>To make the comparison intuitive, for each test context, we dynamically set the beam size of Beam-COMET to the number of facts generated by DIFFUCOMET-Entity.

<sup>&</sup>lt;sup>17</sup>We also include the evaluation results on traditional NLG and our proposed clustering-based metrics in Appendix F.4.

## Limitations

587

589

592

594

595

597

599

601

604

607

610

611

612 613

614

615

616

617

618

619

622

623

624

626

627

629

630 631

632

633

637

We notice a few limitations in this work. First, narrative samples in our experimental datasets, i.e., ComFact (Gao et al., 2022a) and WebNLG+ 2020 (Ferreira et al., 2020), have short context windows (five sentences at maximum). Therefore, our knowledge models trained on these datasets may have limited inference capacities if applied to longer narratives that involve richer commonsense grounding. Moreover, our models are trained on solely English corpora, and may need additional resources to be adapted to other languages or multilingual settings. Finally, our diffusion modeling method is tested on an encoder-decoder model structure, i.e., BART (Lewis et al., 2020), with maximum model size 406M (BART-large). We leave the feasibility of our method on other model structures, e.g. decoderonly GPT (Radford et al., 2019), and larger model scales, to future work.

## References

- David Bamman, Brendan O'Connor, and Noah A Smith.
  2013. Learning latent personas of film characters.
  In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 352–361.
- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
  - Antoine Bosselut, Ronan Le Bras, and Yejin Choi. 2021.
    Dynamic neuro-symbolic knowledge graph construction for zero-shot commonsense question answering.
    In *Proceedings of the AAAI conference on Artificial Intelligence*, volume 35, pages 4923–4931.
  - Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi.
     2019. Comet: Commonsense transformers for automatic knowledge graph construction. In *Proceedings* of the 57th Annual Meeting of the Association for Computational Linguistics, pages 4762–4779.
- Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. 2022. Analog bits: Generating discrete data using diffusion models with self-conditioning. *arXiv preprint arXiv:2208.04202*.
- Leyang Cui, Yu Wu, Shujie Liu, Yue Zhang, and Ming Zhou. 2020. Mutual: A dataset for multi-turn dialogue reasoning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1406–1416.

Jeff Da, Ronan Le Bras, Ximing Lu, Yejin Choi, and Antoine Bosselut. 2021. Analyzing commonsense emergence in few-shot knowledge models. *arXiv preprint arXiv:2101.00297*.

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

681

682

683

684

685

686

687

688

689

690

691

- Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Knowledge discovery and aata mining*, volume 96, pages 226–231.
- Thiago Castro Ferreira, Claire Gardent, Nikolai Ilinykh, Chris Van Der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina. 2020. The 2020 bilingual, bi-directional webnlg+ shared task overview and evaluation results (webnlg+ 2020). In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*.
- Silin Gao, Beatriz Borges, Soyoung Oh, Deniz Bayazit, Saya Kanno, Hiromi Wakaki, Yuki Mitsufuji, and Antoine Bosselut. 2023. Peacok: Persona commonsense knowledge for consistent and engaging narratives. *arXiv preprint arXiv:2305.02364*.
- Silin Gao, Jena D Hwang, Saya Kanno, Hiromi Wakaki, Yuki Mitsufuji, and Antoine Bosselut. 2022a. Comfact: A benchmark for linking contextual commonsense knowledge. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1656–1675.
- Zhujin Gao, Junliang Guo, Xu Tan, Yongxin Zhu, Fang Zhang, Jiang Bian, and Linli Xu. 2022b. Difformer: Empowering diffusion model on embedding space for text generation. *arXiv preprint arXiv:2212.09412*.
- Jian Guan, Yansen Wang, and Minlie Huang. 2019. Story ending generation with incremental encoding and commonsense knowledge. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6473–6480.
- Kehang Han, Kathleen Kenealy, Aditya Barua, Noah Fiedel, and Noah Constant. 2024. Transfer learning for text diffusion models. *arXiv preprint arXiv:2401.17181*.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840– 6851.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*.
- Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. (comet-) atomic 2020: On symbolic and neural commonsense knowledge graphs.

801

802

- In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 6384–6392.
  - Mete Ismayilzada and Antoine Bosselut. 2023. kogito: A commonsense knowledge inference toolkit. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations, pages 96–104.
    - Haozhe Ji, Pei Ke, Shaohan Huang, Furu Wei, Xiaoyan Zhu, and Minlie Huang. 2020. Language generation with multi-hop reasoning on commonsense knowledge graph. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 725–736.

702

710

711

712

713

714

715

716

717

718

719

720

721 722

723

724

725

727

730

731

732

733 734

735

736

737

738

739

740

741

742

743

744

- Liwei Jiang, Antoine Bosselut, Chandra Bhagavatula, and Yejin Choi. 2021. "i'm not mad": Commonsense implications of negation and contradiction. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4380–4397.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. Dbpedia–a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic web*, 6(2):167–195.
- Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals.
  In *Soviet physics doklady*, volume 10, pages 707–710.
  Soviet Union.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.
- Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. 2022. Diffusionlm improves controllable text generation. *Advances in Neural Information Processing Systems*, 35:4328– 4343.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Weizhu Chen, and Nan Duan. 2022. Genie: Large scale pre-training for text generation with diffusion model. *arXiv preprint arXiv:2212.11685*.
- Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

- Igor Melnyk, Pierre Dognin, and Payel Das. 2022. Knowledge graph generation from text. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1610–1622.
- Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849.
- Lassila Ora. 1999. Resource description framework (rdf) model and syntax specification. *http://www.w3.org/TR/REC-rdf-syntax/*.
- 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.
- Xiangyu Peng, Siyan Li, Sarah Wiegreffe, and Mark Riedl. 2022. Inferring the reader: Guiding automated story generation with commonsense reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7008–7029.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- 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. *The Journal of Machine Learning Research*, 21:5485–5551.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for ifthen reasoning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3027–3035.
- Isabel Segura-Bedmar, Paloma Martínez, and María Herrero-Zazo. 2013. Semeval-2013 task 9: Extraction of drug-drug interactions from biomedical texts (ddiextraction 2013). In Second Joint Conference on Lexical and Computational Semantics (\* SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013), pages 341–350.

Jascha Sohl-Dickstein, Eric Weiss, Niru Mah-

eswaranathan, and Surya Ganguli. 2015. Deep un-

supervised learning using nonequilibrium thermo-

dynamics. In International conference on machine

Yang Song and Stefano Ermon. 2019. Generative mod-

eling by estimating gradients of the data distribution.

Advances in neural information processing systems,

Robyn Speer, Joshua Chin, and Catherine Havasi. 2017.

Quan Tu, Yanran Li, Jianwei Cui, Bin Wang, Ji-Rong

Wen, and Rui Yan. 2022. Misc: A mixed strategy-

aware model integrating comet for emotional support conversation. In Proceedings of the 60th Annual

Meeting of the Association for Computational Lin-

guistics (Volume 1: Long Papers), pages 308–319.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz

Peter West, Chandra Bhagavatula, Jack Hessel, Jena

Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu,

Sean Welleck, and Yejin Choi. 2022. Symbolic

knowledge distillation: from general language mod-

els to commonsense models. In Proceedings of the

2022 Conference of the North American Chapter of

the Association for Computational Linguistics: Hu-

man Language Technologies, pages 4602–4625.

Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren,

Xikun Zhang, Christopher D Manning, Percy S Liang, and Jure Leskovec. 2022. Deep bidirectional

language-knowledge graph pretraining. Advances in

Neural Information Processing Systems, 35:37309-

Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut,

Percy Liang, and Jure Leskovec. 2021. Qa-gnn: Reasoning with language models and knowledge graphs

for question answering. In Proceedings of the 2021

Conference of the North American Chapter of the

Association for Computational Linguistics: Human

Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Fei Huang,

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur

Szlam, Douwe Kiela, and Jason Weston. 2018. Per-

sonalizing dialogue agents: I have a dog, do you have

pets too? In Proceedings of the 56th Annual Meet-

ing of the Association for Computational Linguistics

(Volume 1: Long Papers), pages 2204-2213.

and Songfang Huang. 2022. Seqdiffuseq: Text dif-

fusion with encoder-decoder transformers. arXiv

Language Technologies, pages 535–546.

preprint arXiv:2212.10325.

Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing

ence on artificial intelligence, volume 31.

Conceptnet 5.5: An open multilingual graph of gen-

eral knowledge. In Proceedings of the AAAI confer-

learning, pages 2256–2265. PMLR.

- 806 807
- 810 811

812

32.

systems, 30.

37323.

- 813 814
- 816
- 818
- 820 821
- 824
- 827
- 829 830
- 831 832
- 834

838 839

840 841

844

850 851

853 854

- 856 857

Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D. Manning, and Jure Leskovec. 2022. Greaselm: Graph reasoning enhanced language models for question answering. In Proceedings of the 10th International Conference for Learning Representations (ICLR).

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

Pei Zhou, Karthik Gopalakrishnan, Behnam Hedayatnia, Seokhwan Kim, Jay Pujara, Xiang Ren, Yang Liu, and Dilek Hakkani-Tur. 2022. Think before you speak: Explicitly generating implicit commonsense knowledge for response generation. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1237-1252.

#### **Backward Diffusion Process** Α

Inverting from the forward diffusion process formulated as Eq.(1), the backward diffusion process follows a Gaussian posterior distribution  $q(\mathbf{z}_{t-1}|\mathbf{z}_t,\mathbf{z}_0)$ :

$$q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{z}_0) = \mathcal{N}(\mathbf{z}_{t-1}; \widetilde{\mu}(\mathbf{z}_t, \mathbf{z}_0), \widetilde{\beta}_t \mathbf{I})$$
$$\widetilde{\mu}(\mathbf{z}_t, \mathbf{z}_0) = \frac{\sqrt{\overline{\alpha}_{t-1}}\beta_t}{1 - \overline{\alpha}_t} \mathbf{z}_0 + \frac{\sqrt{\alpha_t}(1 - \overline{\alpha}_{t-1})}{1 - \overline{\alpha}_t} \mathbf{z}_t$$
$$\widetilde{\beta}_t = \frac{1 - \overline{\alpha}_{t-1}}{1 - \overline{\alpha}_t} \beta_t$$
(12)

where  $\alpha_t = 1 - \beta_t$  and  $\overline{\alpha}_t = \prod_{i=1}^t \alpha_i$  are weight hyperparameters of the posterior Gaussian defined by the noise schedule  $\beta_t$ . The posterior formulation indicates that only the mean  $\tilde{\mu}$  of  $\mathbf{z}_{t-1}$  is correlated to the condition  $z_t$  and  $z_0$ . So the training loss for diffusion models, derived from the KL-divergence between gold and learned posterior distributions, is typically defined as a mean-squared error loss on the posterior Gaussian mean:

$$\mathcal{L}_{mse} = \sum_{t=1}^{T} \mathbb{E} \| \widetilde{\mu}(\mathbf{z}_t, \mathbf{z}_0) - \mu_{\theta}(\mathbf{z}_t, t) \|^2 \quad (13)$$

where model (with parameter  $\theta$ ) learns the function  $\mu_{\theta}(\mathbf{z}_t, t)$  to predict the mean of  $\mathbf{z}_{t-1}$ . Diffusion-LM (Li et al., 2022) further re-weights the meansquared error as Eq.(2) to enforce direct prediction of  $z_0$  in every loss term, which is shown to be more efficient at tuning the model to precisely predict the final de-noised sample.

#### **Model Implementation Details** B

## **B.1 Diffusion Module**

To conduct the diffusion process defined by Eq.(4) using Transformers (Vaswani et al., 2017),  $z_t$  and

946 947

948

949

950

951

952

953

954

955

956

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

 $\hat{\mathbf{z}}_0^t$  are first concatenated at the hidden-state dimension and projected by a MLP layer to form their joint representation. The positional encoding layer of Transformers is applied to the time step t (same for every position of self-attention), whose output time step embedding is added to the joint representation of  $\mathbf{z}_t$  and  $\hat{\mathbf{z}}_0^t$ . The decoder  $f_{\theta_z}$  takes the joint representation (with time step embedding added) as its bi-directional self-attention input, to ground its decoding of refined  $\mathbf{z}_0$  prediction  $\hat{\mathbf{z}}_0^{t-1}$ .

## B.2 Number of Generated Facts

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

921

922

923

924

926

927

931

932

933

937

939

940

To enable our diffusion module  $(f_{\theta_z})$  to control the number of facts (or entities) generated for each context, we also pre-train our fact (or entity) embedding module  $(f_{\theta_e} \text{ and } f_{\theta_g})$  to learn the representation of a special token  $k_{end} := \langle eok \rangle$ , by adding it as a special fact (or entity) to the pre-training data, which indicates the end of a knowledge set. During the training of diffusion module,  $k_{end}$  is appended to the end of knowledge set  $\mathcal{K}$ , whose embedding and decoding also contributes to the training loss. At inference phase, we post-process our model's generations to keep only the facts that are at positions before  $k_{end}$ .

## B.3 Noise Schedule

For the noise schedule hyperparameter of diffusion process, we adopt the *sqrt* initialization (Li et al., 2022) to set  $\overline{\alpha}_t = 1 - \sqrt{t/T + s}$ , where  $s = 1e^{-4}$ that sets the initial variance of noise ( $\beta_0$ ) to be 0.01. Based on that, we follow SeqDiffuSeq (Yuan et al., 2022) to implement an adaptive noise schedule, which dynamically adjusts  $\overline{\alpha}_t$  for each sample position n ( $n = 1, 2, ... |\mathcal{K}|$ ) of the knowledge set  $\mathcal{K}$ (the adjusted  $\overline{\alpha}_t$  for position n is denoted as  $\overline{\alpha}_t^n$ ), according to the diffusion mean square error (MSE) loss  $\mathcal{L}_{\theta_s, \theta_z}^{mse}$  defined in Eq. (8). Specifically, for an adaptive noise schedule update, we first record the MSE loss at each time t and position n as:

$$\mathcal{L}_{t}^{n} = \mathbb{E} \| \mathbf{z}_{0}[:][n] - \hat{\mathbf{z}}_{0}^{t}[:][n] \|^{2}$$
(14)

Then we use a linear interpolation function to update the adjusted noise schedule, formulated as:

$$F_t^n(\mathcal{L}) = \frac{\overline{\alpha}_t^n - \overline{\alpha}_{t-1}^n}{\mathcal{L}_t^n - \mathcal{L}_{t-1}^n} (\mathcal{L} - \mathcal{L}_{t-1}^n) + \overline{\alpha}_{t-1}^n \quad (15)$$

941 where new loss value  $\mathcal{L}_t^{n,new}$  is re-arranged across 942 time step t with equal interval between  $\min_t(\mathcal{L}_t^n)$ 943 and  $\max_t(\mathcal{L}_t^n)$ , which is finally given to the update 944 function to get  $\overline{\alpha}_t^{n,new} = F_t^n(\mathcal{L}_t^{n,new})$ . The noise 945 schedule is adjusted every 2000 training steps.

#### **B.4 Model Training**

For the loss weight hyperparameter  $\gamma$  used to combine mean-square error and anchor losses defined by Eq. (8) and (9), we use  $\gamma = 1$  for training our DIFFUCOMET models based on BART-base, while  $\gamma = 0.01$  for training our models with BARTlarge backbone, which achieve the best convergence results, respectively. For training DIFFU-COMET based on BART-large, we also follow Difformer (Gao et al., 2022b) to amplify the standard deviation of diffusion noise by a factor of A = 4, *i.e.*, to change the forward process as:

$$q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t A^2 \mathbf{I}) \quad (16)$$

where t = 1, 2, ..., T, which effectively avoids model collapse in training. The total diffusion steps T is set to 2000. We use AdamW (Loshchilov and Hutter, 2018) as our training optimizer, with learning rate  $1e^{-5}$  and no weight decay. A linear learning rate scheduler is adopted with warmup steps 2000 and total training steps 150000 and 300000 for models based on BART-base (139M) and BART-large (406M), respectively. We train our base-scale DIFFUCOMET on 4 Tesla V100-SXM2 (32GB) GPUs with batch size set to 4, while for large-scale DIFFUCOMET, we use 4 NVIDIA A100-SXM4 (40GB) GPUs, with batch size set to 2 instead. 15 and 36 hours are required to train base-scale and large-scale DIFFUCOMET models, respectively.

For the pre-training of our fact embedding module with loss described in Eq. (7), we adopt the same hyperparameter setting as training our diffusion module, except for learning rate changed to  $2e^{-6}$  and batch size set to 128 and 64 for basescale and large-scale models, respectively. For pretraining large-scale (*i.e.*, BART-large) fact embedding module, we add a weight decay of 0.01, which leads to better convergence. In DIFFUCOMET-Entity, the two diffusion modules trained for generating contextual relevant head and tail entities share the same pre-trained entity embedding module.

## **C** Evaluation Metrics

## C.1 Clustering and Similarity Function

For our evaluation based on fact clustering *w.r.t.* edit distance, we define the similarity function in our Alignment metric as  $sim(\hat{k}_i, k_j) = 1 - Edit(\hat{k}_i, k_j)/MaxLen(\hat{k}_i, k_j)$ , where *Edit* denotes the word-level edit distance of two facts, and



Figure 5: Range selection (red square) of DBSCN clustering thresholds for our proposed metrics.

994

995

997

998

999

1002

1005

1006

1007

1008

1010

1011

MaxLen denotes the length of the longer fact of the two, *i.e.*, the maximum possible edit distance for normalization. Our distance measure for clustering also adopts the normalized edit distance, *i.e.*, Edit/MaxLen. For evaluation based on fact clustering w.r.t. Sentence-BERT embedding, we define the similarity function in our Alignment metric as  $sim(k_i, k_j) = max(CoS(k_i, k_j), 0)$ , where CoSdenotes the cosine similarity of two facts' Sentence-BERT embeddings. We assume that facts with opposite meanings, *i.e.*, negative similarity, are not considered as aligned with each other, so we cut off the negative values of cosine similarity. While for the distance measure of clustering, we use the Euclidean distance of two facts' embeddings instead, which is typically adopted in DBSCAN (Ester et al., 1996) clustering algorithm.

## C.2 Clustering Threshold Selection

1012For our proposed clustering-based metrics as de-1013scribed in Section 4, we use DBSCAN (Ester et al.,10141996) algorithm to group facts into clusters. To1015avoid bias on a specific clustering granularity, we1016consider a range of DBSCAN thresholds and take

the average evaluation results across all thresholds 1017 in the range. We consider a range with equal inter-1018 val of 0.05, where the number of gold fact clusters 1019 significantly changes from near the maximum (*i.e.*, 1020 each fact as a cluster) to near the minimum (i.e., 1021 all facts grouped into one cluster). Figure 5 shows 1022 the number of gold clusters as a function of the 1023 DBSCAN clustering threshold, and our selection of threshold ranges (red square) on each dataset. 1025

### C.3 WebNLG Metrics

In the evaluation of WebNLG 2020 Challenge (Fer-1027 reira et al., 2020), each generated RDF fact (i.e., 1028 subject-predicate-object triple) is paired to a gold 1029 reference to compute its precision, recall and F1 1030 based on named entity matching (Segura-Bedmar 1031 et al., 2013). Three types of matching criterias are 1032 considered, including: a) each named entity in gen-1033 erated RDF needs to exactly match an entity in gold 1034 reference in order to be counted as true-positive, 1035 while its type in the RDF (*i.e.*, whether it is in sub-1036 ject, predicate or object) does not need to match 1037 (Exact Match), b) each entity in generated RDF only needs to partially match an entity in gold refer-1039 ence, and its type does not matter (Partial Match), 1040 and c) each named entity in generated RDF needs 1041 to exactly match an entity in gold reference, and 1042 its type also needs to match (Strict Match). For each matching criteria, optimal pairing (with the 1044 highest F1 score) between generated facts and gold 1045 references is searched by enumerating all possible permutations. We report Strict Match scores in the 1047 main body of our paper in Table 4, and include all 1048 three kinds of match scores in Table 18. 1049

## **D** Data Preprocessing

ComFact (Gao et al., 2022a) benchmark con-1051 tains social commonsense knowledge linked from 1052 ATOMIC<sub>20</sub><sup>20</sup> (Hwang et al., 2021) knowledge base, 1053 which contains  $\sim 1.33M$  facts covering physi-1054 cal entities, daily events and social interactions. 1055  $\operatorname{ATOMIC}_{20}^{20}$  commonsense relations considered in our experiments are listed in Table 5. We prepro-1057 cess ComFact and ATOMIC<sup>20</sup><sub>20</sub> to filter out facts 1058 that have invalid tail entity "none" or contain fil-1059 lable blank "\_\_\_\_", *i.e.*, we do not consider facts 1060 with relation "IsFilledBy". After preprocessing, 1061  $\sim$ 972K facts are involved in the training of our fact 1062 embedding and diffusion modules. The original 1063 ComFact training data in the ROCStories portion only has  $\sim 1K$  contexts with gold annotations 1065

Туре	Relation	<b>Relation Description</b>
Physical- Entity	ObjectUse AtLocation MadeUpOf HasProperty CapableOf Desires NotDesires	used for located or found at/in/on made (up) of can be characterized by being/having is/are capable of desires do(es) not desire
Event	IsAfter IsBefore HasSubEvent HinderedBy Causes xReason	happens after happens before includes the event/action can be hindered by causes because
Social- Interaction	xNeed xAttr xEffect xReact xWant xIntent oEffect oReact oWant	but before, person X needs person X is seen as as a result, person X will as a result, person X feels as a result, person X wants because person X wants as a result, others will as a result, others feel as a result, others want

Table 5: Commonsense relations in ATOMIC<sup>20</sup><sub>20</sub> knowledge base that are considered in our experiments on ComFact benchmark.

of relevant facts. Due to the limited supervised data, we augment the training data with ~ 50Kadditional contexts sampled from the ROCStories corpus, and use a DeBERTa (He et al., 2020) fact linker developed from the *ComFact* benchmark to extract silver annotations of relevant facts from ATOMIC<sup>20</sup><sub>20</sub> to each additional context.

1066

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1081

1082

1083

1085

1086

1087

1088

1089

1090

1091

1093

For preprocessing WebNLG+ 2020 (Ferreira et al., 2020) dataset, we follow Grapher (Melnyk et al., 2022) to remove underscores and surrounding quotes appeared in the dataset, and convert non-English characters into their closest available English characters, *e.g.*, "õ" and "å" are mapped to "o" and "a". After preprocessing, We develop our models based on the  $\sim 35K$  WebNLG training texts and their linked RDF facts.

**E** Human Evaluation Details

Our annotator pool for human evaluation contains 58 Amazon Mechanical Turk workers who are located in the USA and have been previously qualified by us for other similar tasks. To prepare the workers for the new tasks of assessing the validity and relevance of knowledge in a given context, we share the instructions with them beforehand and do a small pilot run where we evaluate the quality of the worker annotations and give feedback if needed. We pay each worker \$0.10 for each task. Figure 6, 7 and 8 show screenshots of our acceptance/privacy policy and instructions for knowledge validation1094and relevance tasks. Our data collection protocol1095follows Amazon Mechanical Turk regulations, and1096is approved by our organization in terms of ethics.1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

## F Full Results of Knowledge Generation

## F.1 ROCStories

In Table 6 and 7, we present our full evaluation results of contextual commonsense knowledge generation on the ROCStories portion of *ComFact* benchmark. For evaluating Sampling and Beam baseline models, we test two sampling or beam search sizes around the average number of gold facts per context, *i.e.*, 10 and 15 as indicated by the suffix numbers, and adopt the size which achieves better F1 results. On both base and large model scales, DIFFUCOMET models achieve consistently better balance between the diversity (*i.e.*, # Clusters) and accuracy (*i.e.*, RA-F1) of knowledge generation, compared to baseline models that typically perform generation in the autoregressive manner.

#### F.2 Context Generalization

In this section, we present zero-shot evaluation results of models (trained on the contexts of ROC-Stories) generalizing to the contexts of other three ComFact portions, including PersonaChat (Table 8 and 9), MuTual (Table 10 and 11) and MovieSummaries (Table 12 and 13).

We observe that both DIFFUCOMET models generalize well to the contexts of PersonaChat and MuTual, whose generated knowledge possesses comparable diversity (*i.e.*, # Clusters) and better accuracy (*i.e.*, RA-F1) than the strongest baseline model Beam-COMET. More interestingly, we find that DIFFUCOMET-Entity achieves larger points of improvements over baselines on the more challenging MovieSummaries-style contexts, while DIFFUCOMET-Fact struggles to outperform the strongest baseline Beam-COMET, showing that entity-level diffusion is more robust to the shift of narrative contexts, likely due to the more fine-grained multi-step learning of contextto-knowledge mapping.

#### F.3 Case Study and Knowledge Types

Table 14 showcases the knowledge generation re-<br/>sults of DIFFUCOMET models in a narrative con-<br/>text sampled from ComFact ROCStories, com-<br/>pared to the sampling and beam search baselines11381140

#### Acceptance and Privacy Policies (click to expand/collapse)

#### Acceptance Policy

There is no obligation to participate in the task. We will not reject a job unless we observe the evidence of malicious behavior, such as random clicks or very short session times.

#### Privacy Policy

We may incidentally collect some personal data for the purpose of our research project, according to art. 36c and seq. of the ETH Act. Our target is to process and publish only anonymized data. Raw data will be kept confidential and secure. Only anonymized or aggregated personal data may be shared with other research partners.

Having established this, however, we should not collect any personal data in this task.

We are using the services of Amazon Mechanical Turk, Inc. and its affiliates, c/o Amazon.com, Inc., P.O. Box 81226, Seattle, WA 98108-1226, USA. Hence, the privacy policy of Amazon will apply for the processing of your personal information by Amazon.

If you wish to raise a complaint on how we have handled your personal data, or if you want to know if we hold personal data about you, you can contact our Data Protection Officer (dpo@epfl.ch) who will investigate the matter.

#### Figure 6: Screenshot of Amazon MTurk Acceptance and Privacy Policy

# **Knowledge Validity** Acceptance and Privacy Policies (click to expand/collapse) Instructions (click to expand/collapse) (WARNING: This HIT may contain adult content. Worker discretion is advised.) Thanks for participating in this HIT! Given a list of knowledge statements, you are asked to select statements that are generally possible or valid from the commonsense perspective. An example of list of *knowledge statements* could be the following: I. wrap UsedFor wrap the present 2. gift AtLocation wrapped container 3. ceives gifts xWant to finish wrapping gifts 4. gift UsedFor cooking 5. buys gifts for his family IsBefore wraps gifts 6. paper Not CapableOf wrap gifts ✓ 7. wraps gifts xWant give gifts ✓ 8. buys gifts for his family xEffect wraps the gifts As can be seen from the examples, each knowledge statement is represented as A Relation B where A and B refer to phrases relevant to the context and Relation represents the knowledge relationship between them. We provide list of available knowledge relations below with a brief description and an example below. Given a list of knowledge statements like the above examples, you need to select those ones that are valid in general. What we mean by validness is that the knowledge statement is true or makes sense from our commonsense perspective. Note that we are looking for a soft validity check meaning that you should select a statement if you can somehow interpret it in a sensible way and deselect only if it absolutely does not make sense or it contains a major typo that prevents you from understanding its meaning. In the given example, option 1, 2, 5, 7 and 8 are selected because they are true and sensible statements in general. Option 3 should NOT be selected because it has a major typo (i.e. there is no word "ceives"), so we can't interpret the statement's validity. Option 4 should NOT be selected because it does not generally make sense from the commonsense perspective. Option 6 should NOT be selected because it is not true

Note that phrases are NOT supposed to be complete sentences, but rather entities or short phrases that make sense. Phrases may also contain words like **PersonX** or **PersonY** referring to people in general, so these are NOT typos.

Figure 7: Screenshot of Amazon MTurk instructions for knowledge validation task.

## Knowledge Relevance

Acceptance and Privacy Policies (click to expand/collapse)
Instructions (click to expand/collapse)
(WARNING: This HIT may contain adult content. Worker discretion is advised.) Thanks for participating in this HIT!
Given a short <i>context</i> and a list of <i>knowledge statements</i> , you are asked to <b>select</b> statements that are <b>relevant</b> to the given <i>context</i> . An example of a short <i>context</i> could be the following:
Context hank had to wrap a lot of gifts for his family . he ran out of wrapping paper with 4 gifts to go . he went to the kitchen and found shopping bags . he cut up the bags to make sheets of paper .
An example of list of <i>knowledge statements</i> could be the following:           I.         wrap         UsedFor         wrap         the present
<ul> <li>2. gift AtLocation wrapped container</li> <li>3. ceives gifts xWant to finish wrapping gifts</li> </ul>
4. gift UsedFor make sheets of paper 5. paper CapableOf be published at a conference
<ul> <li>6. buys gifts for his family IsBefore wraps gifts</li> <li>7. buy wrapping paper xIntent package goods for sale</li> </ul>
As can be seen from the examples, each <i>knowledge statement</i> is represented as <b>A Relation B</b> where <b>A</b> and <b>B</b> refer to phrases relevant to the context and <b>Relation</b> represents the knowledge relationship between them. We provide list of available knowledge relations below with a brief description and an example below.
Given a short context and a list of knowledge statements like the above examples, you need to select those ones that are relevant to the context. What we mean by relevance is that the knowledge statement is valid and helpful in understanding the context.
In the given example, option 1, 2 and 6 are selected because they are relevant in understanding the context.
Option 3 should NOT be selected because it is unclear what it means (i.e. there is no word "ceives") and hence is an invalid and irrelevant statement.
<b>Option 4</b> should <b>NOT</b> be selected because it does not make sense and is not true in this context.
Option 5 should NOT be selected because while it is a valid statement in general, it is not really helpful in understanding the context here.

Figure 8: Screenshot of Amazon MTurk instructions for knowledge relevance task.

Backhone	Model	# Facts	Clusterin	g <i>w.r.t</i> . Word	-Level Edit D	istance	Clustering	w.r.t. Embed	ding Euclidea	n Distance
Duckbolic	mouer	" I ucus	# Clusters	Relevance	Alignment	RA-F1	# Clusters	Relevance	Alignment	RA-F1
	Greedy	2.48	1.08	32.04	31.98	32.01	1.09	32.11	48.59	38.67
	Sampling-10	10.00	5.59	39.20	46.03	42.34	5.64	38.93	64.51	48.56
	Sampling-15	15.00	7.64	37.18	49.78	42.57	7.82	36.86	68.00	47.81
BART	Beam-10	10.00	2.63	38.30	44.96	41.36	2.83	38.87	59.58	47.05
(base)	Beam-15	15.00	3.48	41.46	48.04	44.51	3.97	42.88	63.14	51.07
	DIFFUCOMET-Fact	13.40	4.74	59.75	54.07	56.77	5.85	60.32	73.38	66.21
	DIFFUCOMET-Entity	10.08	4.51	62.27	54.61	58.19	5.24	61.77	71.54	66.30
	Greedy	2.20	1.38	60.45	36.11	45.21	1.37	60.22	52.31	55.99
	Sampling-10	10.00	6.68	56.09	52.10	54.02	6.40	56.68	73.86	64.14
DADT	Sampling-15	15.00	8.89	56.24	55.18	55.70	8.56	56.57	76.30	64.97
(larga)	Beam-10	10.00	3.32	64.94	50.72	56.96	3.51	64.37	69.14	66.67
(large)	Beam-15	15.00	4.17	64.18	53.66	58.45	4.60	64.35	71.35	67.67
	DIFFUCOMET-Fact	12.88	4.47	65.82	54.18	59.44	5.24	65.64	71.65	68.51
	DIFFUCOMET-Entity	12.89	5.09	67.00	58.22	62.30	5.67	66.39	74.38	70.16
	Greedy	1.96	1.14	61.27	34.76	44.36	1.19	61.42	50.64	55.51
CONT	Sampling-10	10.00	6.45	56.79	53.36	55.02	6.30	56.60	73.64	64.01
DADT	Sampling-15	15.00	8.52	55.78	58.99	57.34	8.39	56.19	77.97	65.31
DAKI	Beam-10	10.00	3.78	65.62	53.45	58.91	3.89	65.73	70.65	68.10
	Beam-15	15.00	4.78	64.91	54.77	59.41	5.09	65.03	71.64	68.18
T5 (large)	Grapher	5.08	1.75	67.82	33.07	44.46	2.60	68.29	40.58	50.91
-	Gold	10.55	5.64	81.06	-	-	5.64	80.90	-	-

Table 6: Clustering-based evaluation results on the **ROCStories** portion of ComFact. Best results (excluding Gold references) are in bold. Different numbers after Sampling and Beam denote various sampling numbers or beam search sizes being tested.

Backbone	Model	Distinct-4	BLEU	METEOR	ROUGE-L
	Greedy	99.90	8.70	40.49	44.43
	Sampling-10	85.29	7.16	37.78	39.20
DADT	Sampling-15	81.57	8.24	38.35	40.13
(hasa)	Beam-10	50.32	12.25	42.23	43.53
(Dase)	Beam-15	45.21	11.51	42.04	42.91
	DIFFUCOMET-Fact	57.87	12.09	46.43	47.13
	DIFFUCOMET-Entity	70.02	14.25	43.34	45.08
	Greedy	93.01	9.12	43.98	46.26
	Sampling-10	86.33	9.89	43.85	43.69
DADT	Sampling-15	81.56	9.47	43.28	43.15
BARI	Beam-10	47.03	15.02	48.56	48.15
(large)	Beam-15	43.73	13.11	47.70	46.35
	DIFFUCOMET-Fact	52.46	15.98	50.06	51.44
	DIFFUCOMET-Entity	63.49	17.01	47.61	48.40
	Greedy	65.95	18.01	52.32	54.96
Court	Sampling-10	83.29	13.35	44.77	45.80
DADT	Sampling-15	79.01	12.69	44.43	45.58
BARI	Beam-10	51.13	19.89	50.14	50.48
	Beam-15	47.27	16.97	47.39	47.19
T5 (large)	Grapher	67.83	1.40	23.96	27.21
-	Gold	80.45	-	-	-

Table 7: Evaluation results of natural language generation metrics on the **ROCStories** portion of ComFact. Notations are same as Table 6.

Sample-COMET and Beam-COMET. Facts that are novel (*i.e.*, beyond the coverage of gold references) and relevant to the context are labeled in bold. We find that both DIFFUCOMET-Fact and DIFFUCOMET-Entity can generate facts that are rich in diversity, covering both physical entities (*e.g.*, baseball cap) and social events (*e.g.*, go on vacation). Novel facts generated by DIFFUCOMET models also uncover implicit inter-connections between entities or events in the narrative context,

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

*e.g.*, "vacation" and "family" are associated because "X goes on vacation" to "spend time with family". By contrast, Beam-COMET model mainly generates simple facts about physical entities, and Sample-COMET model generates many facts that are irrelevant to the context, *e.g.*, "field is used for playing baseball".

1152

1153

1154

1155

1156

1157

1158

We also conduct a study on the proportion of 1159 different knowledge types that each model gener-1160 ates per context, based on the ROCStories portion 1161 of ComFact benchmark. In particular, we divide 1162 commonsense facts into three types according to 1163 their relation groups under  $\text{ATOMIC}_{20}^{20}$  knowledge 1164 scheme, as shown in Table 5, including facts that 1165 are centered on physical entities, events and social 1166 interactions. Table 15 shows the results of knowl-1167 edge proportion generated by DIFFUCOMET and 1168 baseline models, with gold references. Compared 1169 to Sampling-COMET and Beam-COMET baselines, 1170 DIFFUCOMET models generate a larger propor-1171 tion of facts that reveal complex event or social 1172 inter-connections. The proportion of social-based 1173 facts generated by DIFFUCOMET even signifi-1174 cantly surpasses the gold references. All above re-1175 sults imply that diffusion models have the potential 1176 to uncover more in-depth and implicit common-1177 sense inferences from narrative contexts, which 1178 may not be easily extracted from existing knowl-1179 edge bases. 1180

Backbone	Model	# Facts	Clustering	g <i>w.r.t</i> . Word	-Level Edit D	istance	Clustering	w.r.t. Embed	ding Euclidea	n Distance
Duchoone			# Clusters	Relevance	Alignment	RA-F1	# Clusters	Relevance	Alignment	RA-F1
	Greedy	2.62	1.24	33.52	35.64	34.55	1.24	33.57	48.61	39.71
DADT	Sampling-10	10.00	5.23	29.35	44.49	35.37	4.63	28.26	58.85	38.18
DAKI	Beam-15	15.00	3.65	31.11	47.01	37.44	3.49	30.41	58.90	40.11
(base)	DIFFUCOMET-Fact	13.73	4.97	44.25	52.81	48.15	5.24	44.76	69.24	54.37
	DIFFUCOMET-Entity	11.40	4.99	50.39	54.84	52.52	4.94	49.36	68.55	57.39
DADT	Beam-15	15.00	4.44	53.98	54.13	54.05	4.04	54.07	63.17	58.27
DARI	DIFFUCOMET-Fact	10.82	4.48	55.44	55.20	55.32	3.89	55.02	65.20	59.68
(large)	DIFFUCOMET-Entity	12.06	4.72	55.08	57.12	56.08	4.48	54.42	68.11	60.50
COMET-BART	Beam-15	15.00	4.86	54.02	54.78	54.40	4.27	53.97	65.15	59.04
T5 (large)	Grapher	4.53	1.68	47.74	30.51	37.23	1.57	47.94	36.18	41.24
-	Gold	8.60	4.76	70.42	-	-	4.28	70.42	-	-

Table 8: Zero-shot clustering-based evaluation results on the **PersonaChat** portion of ComFact. Notations are same as Table 6.

Backbone	Model	Distinct-4	BLEU	METEOR	ROUGE-L
	Greedy	97.83	8.72	44.44	46.52
DADT	Sampling-10	86.81	4.09	32.95	33.80
(hear)	Beam-15	53.05	8.06	37.43	38.62
(base)	DIFFUCOMET-Fact	63.40	5.84	37.53	39.33
	DIFFUCOMET-Entity	73.38	9.04	34.35	36.46
DADT	Beam-15	47.64	8.71	41.40	40.44
(larga)	DIFFUCOMET-Fact	57.23	8.05	45.83	47.11
(large)	DIFFUCOMET-Entity	68.54	11.11	38.88	40.04
COMET-BART	Beam-15	50.13	10.25	43.47	42.38
T5 (large)	Grapher	52.99	0.68	19.91	22.41
-	Gold	84.96	-	-	-

Table 9: Zero-shot evaluation results of natural language generation metrics on the **PersonaChat** portion of ComFact. Notations are same as Table 6.

## F.4 WebNLG+ 2020

1181

1182

1183 1184

1185

1186

1187

1188 1189

1190

1191

1192

1193

1194

We present our full evaluation results on the WebNLG+ 2020 benchmark in Table 16, 17 and 18. For evaluating Sampling and Beam baselines, we set both sampling and beam search sizes as 5, which is around the average number of gold facts per context. Consistent with the evaluation results on ComFact, DIFFUCOMET models keep achieving better performances on the WebNLG task of factual knowledge generation, implying that our method of diffusion-based contextual knowledge generation can generalize well to knowledge beyond commonsense.

## F.5 Comparison of Fact and Entity Diffusion

For the comparison in between our two diffusion 1195 models, DIFFUCOMET-Entity in general outper-1196 forms DIFFUCOMET-Fact on our proposed met-1197 rics, which may benefit from more fine-grained 1198 1199 multi-step learning of knowledge construction in pipeline. However, DIFFUCOMET-Fact is compu-1200 tational cheaper, *i.e.*, only requires a single step of 1201 fact diffusion instead of two steps of (head and tail) entity diffusion and a relation prediction. 1203

## G Claim of Usage

Our use of existing scientific artifacts cited in this1205paper is consistent with their intended use. Our1206developed code and models are intended to be used1207for only research purposes, any usage of our sci-1208entific artifacts that is outside of research contexts1209should not be allowed.1210

Backbone	Backbone Model		Clustering w.r.t. Word-Level Edit Distance				Clustering w.r.t. Embedding Euclidean Distance			
Duchoone			# Clusters	Relevance	Alignment	RA-F1	# Clusters	Relevance	Alignment	RA-F1
	Greedy	2.50	1.16	35.85	33.63	34.70	1.19	36.10	48.53	41.40
DADT	Sampling-10	10.00	5.68	43.77	45.76	44.74	5.88	43.98	63.75	52.05
DAKI	Beam-15	15.00	3.55	41.27	49.72	45.10	4.08	42.27	63.17	50.65
(base)	DIFFUCOMET-Fact	13.11	4.54	57.64	52.25	54.81	5.57	57.51	70.10	63.18
	DIFFUCOMET-Entity	10.63	4.60	60.08	54.65	57.24	5.27	59.08	68.88	63.60
DADT	Beam-15	15.00	3.92	64.19	51.75	57.30	4.45	62.41	67.31	64.77
DAKI (larga)	DIFFUCOMET-Fact	10.46	4.33	64.74	54.51	59.19	4.80	64.13	68.07	66.04
(large)	DIFFUCOMET-Entity	11.85	4.70	64.39	55.91	59.85	5.39	63.82	71.22	67.32
COMET-BART	Beam-15	15.00	4.52	61.88	54.04	57.69	4.75	60.56	69.72	64.82
T5 (large)	Grapher	4.50	1.70	73.30	32.74	45.26	1.78	73.33	43.13	54.31
-	Gold	10.80	5.58	74.63	-	-	5.79	74.77	-	-

Table 10: Zero-shot clustering-based evaluation results on the **MuTual** portion of ComFact. Notations are same as Table 6.

Backbone	Model	Distinct-4	BLEU	METEOR	ROUGE-L
	Greedy	97.47	14.05	49.89	50.78
DADT	Sampling-10	86.42	5.61	37.11	38.60
(hasa)	Beam-15	49.31	11.57	45.47	45.51
(base)	DIFFUCOMET-Fact	60.66	8.71	44.23	46.11
	DIFFUCOMET-Entity	70.94	11.08	40.28	42.15
DADT	Beam-15	43.91	11.37	46.86	46.75
BAR1 (large)	DIFFUCOMET-Fact	52.00	12.33	49.50	50.97
(large)	DIFFUCOMET-Entity	66.12	12.68	45.11	45.73
COMET-BART	Beam-15	47.40	12.40	49.12	48.57
T5 (large)	Grapher	51.30	1.96	24.70	29.36
-	Gold	80.99	-	-	-

Table 11: Zero-shot evaluation results of natural language generation metrics on the **MuTual** portion of ComFact. Notations are same as Table 6.

Backbone	# Facts	Clustering	g <i>w.r.t</i> . Word	-Level Edit D	istance	Clustering w.r.t. Embedding Euclidean Distance				
Duchoone			# Clusters	Relevance	Alignment	RA-F1	# Clusters	Relevance	Alignment	RA-F1
	Greedy	2.59	1.12	33.50	25.28	28.82	1.11	33.38	37.95	35.52
DADT	Sampling-10	10.00	4.90	26.61	33.31	29.59	4.24	24.96	50.26	33.36
DAKI	Beam-15	15.00	3.54	30.45	36.07	33.02	3.17	29.13	49.63	36.71
(base)	DIFFUCOMET-Fact	14.61	6.02	35.97	38.76	37.31	5.49	36.29	59.83	45.18
	DIFFUCOMET-Entity	15.82	6.31	39.86	39.93	39.89	5.76	39.57	57.55	46.90
DADT	Beam-15	15.00	4.46	42.92	32.79	37.18	3.70	42.52	50.46	46.15
DAKI (larga)	DIFFUCOMET-Fact	8.29	3.01	41.50	30.47	35.14	2.89	40.82	46.59	43.51
(large)	DIFFUCOMET-Entity	13.50	6.28	44.08	40.46	42.19	5.84	42.70	61.56	50.42
COMET-BART	Beam-15	15.00	5.06	41.97	34.63	37.95	4.04	41.54	51.24	45.88
T5 (large)	Grapher	5.34	1.83	54.09	23.74	33.00	1.50	54.12	34.27	41.97
-	Gold	9.00	5.64	58.55	-	-	4.81	58.37	-	-

Table 12: Zero-shot clustering-based evaluation results on the **MovieSummaries** portion of ComFact. Notations are same as Table 6.

Backbone	Model	Distinct-4	BLEU	METEOR	ROUGE-L
	Greedy	95.18	5.14	34.24	36.53
DADT	Sampling-10	90.99	2.52	24.56	28.08
(base)	Beam-15	53.65	4.82	28.33	31.56
(base)	DIFFUCOMET-Fact	63.93	2.27	27.09	30.08
	DIFFUCOMET-Entity	67.24	2.68	24.14	26.95
DADT	Beam-15	47.57	4.89	31.41	33.19
(large)	DIFFUCOMET-Fact	43.68	5.26	34.55	38.80
(large)	DIFFUCOMET-Entity	67.13	3.83	26.63	29.32
COMET-BART	Beam-15	50.29	5.18	31.36	33.39
T5 (large)	Grapher	42.76	0.46	18.24	21.28
-	Gold	87.39	-	-	-

Table 13: Zero-shot evaluation results of natural language generation metrics on the **MovieSummaries** portion of ComFact. Notations are same as Table 6.

Narrative Context	Dustin loved to wear his baseball cap everywhere he went. On vacation his family visited the windy city of Chicago. Dustin's baseball cap blew off his head and into the street. His dad waited until it was safe before getting Dustin's cap. He loved his baseball cap even though it was a little dirty.
Gold	cap, used for, to wear on head cap, used for, wear on their heads head cap, used for, put on head vacation, used for, have fun on vacation, used for, fun vacation, used for, relax out of work and school family, is capable of, plan to go on vacation X takes a family trip, because X wants, to go on vacation X visits the city, X is seen as, traveling dad, can be characterized by being, one of human's parents
Sample -COMET	baseball cap, used for, protect the head baseball cap, used for, protect your head while playing baseball baseball cap, used for, wearing over head X ops for baseball, but before X needs, to find a baseball X's favorite baseball, but before X needs, to find a baseball X's favorite baseball, because X wants to, enjoy the sport baseball, used for, sport as a mascot vacation, used for, have fun on X chases the wind, because X wants, to walk around port, used for, playing baseball cap, used for, playing baseball cap, used for, protection from wind cap, used for, protect head while traveling cap, used for, wear around head jersey, used for, wear while playing
Beam -COMET	baseball cap, used for, wear while playing baseball baseball cap, used for, wear on their head baseball cap, used for, wear on the head baseball, used for, playing baseball with friends baseball, used for, playing baseball with family sport cap, used for, wear while playing <b>Chicago, can be characterized by having, many streets</b> <b>Chicago, can be characterized by having, many streets</b> <b>Chicago, can be characterized by having, many neighborhoods</b> cap, used for, wear on head while playing baseball cap, used for, protect head from wind cap, used for, protect head from wind blows <b>cap, used for, keep the cap on</b> <b>cap, used for, keep the cap clean</b>
DIFFUCOMET -Fact	baseball cap, used for, to put on baseball cap, used for, to keep baseball cap on head baseball cap, used for, wear baseball cap, used for, to play baseball with city, used for, live in vacation, used for, relax X takes a family trip, but before X needs, to spend time with family X takes a family trip, because X wants, to enjoy family time X goes on vacation, because X wants, to spend time with family dad, can be characterized by being, one of human's parents dad's car, used for, to be safe safe, used for, safe to wear
DIFFUCOMET -Entity	cap, used for, wear on head cap, used for, wear on the head baseball cap, used for, look professional baseball cap, used for, to play baseball with X is wearing cap, but before X needs, have a cap X is wearing cap, but before X needs, put on a cap go on vacation, includes the action, take family to beach go on vacation, includes the action, go somewhere nice vacation, used for, enjoy your time off X goes on vacation, because X wants, to spend time with family dad, can be characterized by being, one of human's parents safe, used for, keeping things safe

Table 14: Examples of contextual knowledge generation. Novel and contextually relevant facts are in bold. Model notations are same as Table 1.

Model	Physical	Event	Social
Sampling-COMET	46.17	4.01	49.82
Beam-COMET	60.72	2.36	36.92
DIFFUCOMET-Fact	41.00	4.66	54.34
DIFFUCOMET-Entity	35.75	4.53	59.72
Gold	43.54	7.32	49.14

Table 15: Proportion (%) of different types of knowledge generation on the ROCStories portion of ComFact. "Physical", "Event" and "Social" denote facts with relation types belonging to physical-entity, event and social-interaction, respectively, as shown in Table 5. Model notations are same as Table 1.

Backbone	Model	# Facts	Clusterin	g <i>w.r.t</i> . Word	-Level Edit D	istance	Clustering w.r.t. Embedding Euclidean I			n Distance
Duchoone			# Clusters	Relevance	Alignment	RA-F1	# Clusters	Relevance	Alignment	RA-F1
	Greedy	1.69	0.88	83.16	50.26	62.65	0.88	83.16	71.71	77.01
DADT	Sampling-5	5.00	2.09	81.10	71.25	75.86	1.73	80.89	83.93	82.38
(larga)	Beam-5	5.00	2.12	82.70	72.69	77.37	1.64	82.50	85.78	84.11
(large)	DIFFUCOMET-Fact	2.56	1.69	84.39	74.12	78.92	1.51	84.38	86.18	85.27
	DIFFUCOMET-Entity	2.71	1.82	87.86	78.46	82.89	1.57	87.76	88.59	88.17
CONT	Greedy	1.61	0.96	83.33	54.34	65.78	0.95	83.33	77.23	80.16
DADT	Sampling-5	5.00	2.09	80.89	72.77	76.62	1.76	80.72	84.21	82.43
BARI	Beam-5	5.00	2.15	82.11	72.94	77.25	1.70	81.94	85.94	83.89
T5 (large)	Grapher	2.10	1.39	83.48	70.66	76.54	1.29	83.46	82.21	82.83
-	Gold	3.22	2.27	96.43	-	-	1.91	96.43	-	-

Table 16: Clustering-based evaluation results on the WebNLG+ 2020 benchmark. Notations are same as Table 6.

Backbone	Model	Distinct-4	BLEU	METEOR	ROUGE-L
	Greedy	87.29	81.12	84.57	84.92
DADT	Sampling-5	48.24	74.22	81.71	81.19
(larga)	Beam-5	45.58	75.01	81.78	80.51
(large)	DIFFUCOMET-Fact	81.02	80.43	83.23	84.30
	DIFFUCOMET-Entity	82.20	83.04	89.88	89.72
CONTE	Greedy	93.17	81.43	84.95	85.34
BART	Sampling-5	47.36	75.44	81.84	81.85
	Beam-5	46.17	73.56	80.93	79.46
T5 (large)	Grapher	89.95	76.17	79.61	80.89
-	Gold	82.05	-	-	-

Table 17: Evaluation results of natural language generation metrics on the **WebNLG+ 2020** benchmark. Notations are same as Table 6.

Backbone	Model	E	xact Match		Pa	Partial Match Strict Matc			trict Match	
	Widden	Web-Prec.	Web-Rec.	Web-F1	Web-Prec.	Web-Rec.	Web-F1	Web-Prec.	Web-Rec.	Web-F1
	Greedy	50.42	52.79	51.51	53.76	56.84	55.20	50.14	52.53	51.25
DADT	Sampling-5	73.65	76.73	75.11	79.57	83.89	81.66	72.37	75.45	73.83
BARI	Beam-5	75.32	78.39	76.76	81.32	85.72	83.38	73.36	76.27	74.75
(large)	DIFFUCOMET-Fact	76.59	78.35	77.47	79.17	81.52	80.35	76.30	78.07	77.19
	DIFFUCOMET-Entity	80.80	82.97	81.84	83.72	86.48	85.07	80.68	82.89	81.74
Covern	Greedy	52.55	54.82	53.62	55.99	58.95	57.39	52.30	54.59	53.37
DADT	Sampling-5	74.96	77.87	76.33	80.31	84.41	82.18	73.77	76.67	75.15
BARI	Beam-5	75.95	78.88	77.03	81.66	85.84	83.15	73.80	76.61	74.85
T5 (large)	Grapher	71.50	73.30	72.20	74.10	76.50	75.00	71.20	73.00	71.90

Table 18: Evaluation results on official metrics provided by the **WebNLG+ 2020** benchmark challenge. We present the results of Grapher as reported in its paper. Notations are same as Table 6.