

DoCoGen: Domain Counterfactual Generation for Low Resource Domain Adaptation

Anonymous ACL submission

Abstract

Natural language processing (NLP) algorithms have become very successful, but they still struggle when applied to out-of-distribution examples. In this paper we propose a controllable generation approach in order to deal with this domain adaptation (DA) challenge. Given an input text example, our DoCoGen algorithm generates a domain-counterfactual textual example (D-CON) – that is similar to the original in all aspects, including the task label, but its domain is changed to a desired one. Importantly, DoCoGen is trained using only unlabeled examples from multiple domains – no NLP task labels or pairs of textual examples and their domain-counterfactuals are required. We use the D-CONS generated by DoCoGen to augment a sentiment classifier in 20 DA setups, where source-domain labeled data is scarce. Our model outperforms strong baselines and improves the accuracy of a state-of-the-art unsupervised DA algorithm.¹

1 Introduction

Natural Language Processing (NLP) algorithms are constantly improving and reaching significant milestones (Devlin et al., 2019; Raffel et al., 2020; Brown et al., 2020). However, such algorithms rely on the availability of sufficient labeled data and the assumption that the training and test sets are drawn from the same underlying distribution. Unfortunately, these assumptions do not hold in many cases due to the costly and labor-intensive data labeling process and since text may originate from many different domains. As generalization in low resource regimes and beyond the training distribution are still fundamental NLP challenges, NLP algorithms significantly degrade when applied to such scenarios.

Domain adaptation (DA) is an established field of research in NLP (Roark and Bacchiani, 2003;

Daumé III and Marcu, 2006; Reichart and Rappoport, 2007) that attempts to explicitly address generalization beyond the training distribution (§2). DA algorithms are trained on annotated data from source domains to be effectively applied in various target domains. Indeed, DA algorithms have been developed for multiple NLP tasks throughout the last two decades (Blitzer et al., 2006, 2007; Glorot et al., 2011; Rush et al., 2012; Ziser and Reichart, 2017, 2018a,b; Han and Eisenstein, 2019).

A natural alternative to costly human annotation would be to automatically generate labeled examples for model training. Doing so may expose the model to additional training examples and better represent the data distribution within and outside the annotated source domains. Unfortunately, generating labeled textual data is challenging (Feng et al., 2021), especially when the available labeled data is scarce. Indeed, labeled data generation has hardly been applied to DA (§2).

To allow DA through labeled data generation, we present DoCoGen, an algorithm that generates domain-counterfactual textual examples (D-CONS). In order to do that, DoCoGen intervenes on the domain-specific terms of its input example, replacing them with terms that are relevant for its target domain while keeping all other properties fixed, including the task label. Consider the task of sentiment classification (top example in Table 1). When DoCoGen encounters an example from the *Kitchen* domain (its source domain), it first recognizes the terms related to *Kitchen* reviews, i.e., *knife* and *solid*. Then, it intervenes on these terms, replacing them with text that connects the example to the *Electronics* domain (its target domain) while keeping the negative sentiment.

DoCoGen is a *controllable generation* algorithm (Li et al., 2016; Russo et al., 2020) that is trained using a novel *unsupervised* sentence reconstruction objective. Importantly, it does not require task-annotated data, or pairs of sentences and their

¹Our code and data will be released upon acceptance.

D-CONS. A key component of DoCoGen is the *domain orientation vector*, which guides the model to generate the new text in the desired domain. The parameters of the orientation vectors are learned during the unsupervised training process, allowing the generation model to share information among the various domains it is exposed to.

We focus on two low resource scenarios: Unsupervised domain adaptation (UDA) and any domain adaptation (ADA, Ben-David et al. (2021)), with only a handful of labeled examples available from a single source domain. In both UDA and ADA the model is exposed to limited labeled source domain data and to unlabeled data from several domains. However, in UDA the *unlabeled domains* contain the future target domain to which the model will be applied, while in ADA the model has no access to the target domain during training. To cope with these extreme conditions, we use DoCoGen to enrich the source labeled data with D-CONS from the unlabeled domains. By introducing labeled D-CONS from various domains, we hope to provide the model with a training signal that is less affected by spurious correlations: Correlations between features and the task label which do not hold out-of-domain (OOD) (Veitch et al., 2021).

After a brief evaluation of the intrinsic quality of the D-CONS generated by DoCoGen, we evaluate our complete DA pipeline. We focus on sentiment classification in the UDA and ADA scenarios, for a total of 12 UDA and 8 ADA setups. Our results demonstrate the superiority of DoCoGen over strong DA and textual-data augmentation algorithms. Finally, combining DoCoGen with PERL (Ben-David et al., 2020), a SOTA UDA model, yields new SOTA DA accuracy and stability.

2 Related Work

We first describe research in our DA setups: UDA and ADA. We then continue with the study of counterfactual-based data augmentation, and, finally, we describe research on counterfactual generation methods.

Domain Adaptation (DA) The NLP literature contains several DA setups, the most realistic of which is *unsupervised domain adaptation* (UDA), which assumes the availability of unlabeled data from a source and a target domain, as well as access to labeled data from the source domain (Blitzer et al., 2006). An even more challenging and potentially more realistic setup is the recently proposed

any domain adaptation setup (ADA, Ben-David et al. (2021)), which assumes no knowledge of the target domains at training time. There are several approaches to DA, including representation learning (Blitzer et al., 2006; Ziser and Reichart, 2017) and data-centric approaches like instance re-weighting and self-training (Huang et al., 2006; Rotman and Reichart, 2019).

Since the rise of deep neural networks (DNNs), most focus in DA research has been directed to deep representation learning approaches (DReL). One line of DReL work employs an input reconstruction objective (Glorot et al., 2011; Chen et al., 2012; Yang and Eisenstein, 2014; Ganin et al., 2016). Another line employs pivot features, which are prominent to the task of interest and common in the source and target domains (Blitzer et al., 2007; Pan et al., 2010; Ziser and Reichart, 2018b; Ben-David et al., 2020; Lekhtman et al., 2021).

We deviate from the DReL approach to DA and propose a data-centric methodology. Contrary to the above works, our approach can be applied to both UDA and ADA. Moreover, unlike previous ADA work, which builds upon multi-source DA, our approach can also perform single-source ADA.

Counterfactually Augmented Data (CAD)

Textual data augmentation (TDA) is a technique for increasing the training dataset without explicitly collecting new examples. This is achieved by adding slightly modified copies of already existing examples (local sampling) or newly created data (global sampling). TDA serves as a solution for insufficient data scenarios and as a technique for improving model robustness (Xie et al., 2020; Ng et al., 2020). There are rule-based and model-based approaches to TDA. Rule-based methods commonly involve insertion, deletion, swap and replacement of specific words (Wei and Zou, 2019), or template-based paraphrasing (Rosenberg et al., 2021). Model-based methods typically utilize a pretrained language model (PLM), e.g., for replacing random words (Kobayashi, 2018; Ng et al., 2020), or generating entirely new examples from a prior data-distribution (Bowman et al., 2016; Russo et al., 2020; Wang et al., 2021). Other model-based methods apply backtranslation (Edunov et al., 2018) or paraphrasing (Kumar et al., 2019) for local sampling.

Another approach within local sampling TDA is to change (only) a specific concept that exists in the original example, creating a counterfactual ex-

Original, **Kitchen**: A good **knife** but Quality Control was poor. The **knife** is **solid** and very comfortable in hand, however, when I got it new, the **blade** is slightly **bent**. I expect it to be in almost Perfect **condition**, but it's not.

DoCoGen, **Kitchen** → **Electronics**: A good **product** but Quality Control was poor. The **ipod** is **very easy to use** and very comfortable in hand, however, when I got it new, the **ipod** is slightly **flimsy**. I expect it to be in almost perfect **shape**, but it's not.

Original, **DVD**: The **direction** of this **film** is excellent. I love **all the characters** and the way they interact. The **storyline** is very important also. It's **about religious beliefs** and neighbors that **interact with** each other. It's a well-**paced** and **interesting story** that's not like anything else I've ever **seen**.

DoCoGen, **DVD** → **Airline**: The **service on this flight** is excellent. I love **the staff** and the way they interact. The **safety** is very important also. It's **nice to have staff** and neighbors that **can help** each other. It's a well-**groomed** and **professional crew** that's not like anything else I've ever **experienced**.

Table 1: Domain-counterfactual textual examples (D-CONS) generated by DoCoGen. Red terms are replaced with green terms through the process of D-CON generation. For additional examples see §A.

ample. Counterfactually-Augmented Data (CAD) is generated by minimally intervening on examples to change their ground-truth label, that is, perturbing only those terms necessary to change the label (Kaushik et al., 2020). CAD is commonly used to improve generalizability (Kaushik et al., 2020; Sen et al., 2021), however empirical results using CAD for OOD generalization have been mixed (Joshi and He, 2021; Khashabi et al., 2020).

In this work, we explore a different type of counterfactuals, namely D-CONS, which are the result of intervening only on the example’s domain while holding everything else equal, particularly its task label. For sentiment analysis, we may be, for example, interested in revising a negative movie review, making it a negative airline review. In addition, while CAD is mostly generated via a human-in-the-loop process (Kaushik et al., 2020; Khashabi et al., 2020; Sen et al., 2021), our work focuses on automatic counterfactual generation.

Counterfactual Generation *controllable generation* refers to generation of text while controlling for specific attributes (Prabhumoye et al., 2020). The controlled attributes can range from style (e.g., politeness and sentiment) to content (e.g., keywords and entities) and even topic. Keskar et al. (2019) propose to control the generated text by training an LM on datasets annotated with the controlled attributes, and Meister et al. (2020) modify the model’s decoding method. Recently, Russo et al. (2020) introduced a global sampling conditional variational autoencoder (VAE), augmenting text while controlling for attributes such as label and verb tense. However, controlling for the task label is challenging in scarce labeled data scenarios (Chen et al., 2021), since generative models require large amounts of labeled data .

Counterfactual generation lies at the intersection of controllable generation and causal inference (Feder et al., 2021a). Only few works deal with counterfactual generation, mostly by intervening on the task label. Wu et al. (2021) train a model on textual examples and their manually generated counterfactuals. Other works present methods for controlling for the text domain and semantics (Wang et al., 2020; Feng et al., 2019), yet they all experiment with short texts. A recent work by Yu et al. (2021) focuses on generation of new target-domain examples for aspect-based sentiment analysis (ABSA) (Pontiki et al., 2016). However, this method is designed specifically for ABSA, utilizing predefined knowledge, and is only suitable for UDA setups where source domain labeled data is abundant. Our work presents a novel domain counterfactual generation algorithm, which can be trained in an unsupervised manner, and its generated outputs are demonstrated to be effective in low-resource DA scenarios.

3 Domain-Counterfactual Examples

In this section, we formally define the concept of domain-counterfactual textual examples (D-CONS) and discuss the motivation behind them.

Definition x' is a *domain-counterfactual example* (D-CON) of x if it is a coherent human-like text that is a result of intervening on the domain of x and changing it to another domain, while holding everything else equal. Particularly, we would like the task label of x' and x to be identical. Formally, given an example $(x, y) \sim \mathcal{D}$ and a destination domain \mathcal{D}' , the goal of D-CON generation is to generate $x' \sim P_{\mathcal{D}'}(X|Y = y)$ such that $x' \simeq_{\mathcal{D}'} x$, where $\simeq_{\mathcal{D}'}$ is the domain counterfactual operator.

In this work, given a labeled source example x

we aim to generate coherent human-like D-CONS from the unlabeled domains (see §1). We propose a D-CON generation algorithm, DoCoGen, consisting of two components. The first involves masking domain specific terms of the given example, yielding $M(x)$. The second is a controllable generation model G which takes as input $M(x)$ and a *domain orientation vector* v' . This vector specifies the destination domain \mathcal{D}' , controlling the semantics of the generated D-CON. Formally:

$$\text{DoCoGen}(x, \mathcal{D}') = G(M(x), v') \simeq_{\mathcal{D}'} x$$

Motivation The NLP community has recently become increasingly concerned with *spurious correlations* (Geirhos et al., 2020; Wang and Culotta, 2020; Gardner et al., 2021). In the case of DA, spurious correlations may be defined as correlations between X and Y which are relevant only to a specific domain or in a certain sample of labeled examples. Such correlations may make a predictor $f : \mathcal{X} \rightarrow \mathcal{Y}$ brittle to domain shifts.

Using counterfactuals w.r.t. a specific variable allows us to both estimate its effect on our predictor (Feder et al., 2021b; Rosenberg et al., 2021) or alleviate its impact on it (Kaushik et al., 2021). We focus on the latter, automatically generating D-CONS by intervening on the domain variable \mathcal{D} . Adding these D-CONS to the training set of a predictor should reduce its reliance on domain-specific information and spurious correlations.

From a DA perspective, enriching the training data with D-CONS is motivated by pivot features (§2), which are frequent in multiple domains and are prominent for the task. D-CONS preserve language patterns, such as pivots, which are frequent in multiple domains. Consider the bottom example in Table 1, pivot words (such as *excellent* and *important*) are preserved in the D-CON, while non-pivots (*interesting* and *well-paced*) are replaced due to the domain intervention. Accordingly, a model trained on an example and its D-CON is directed to focus on pivots rather than on non-pivots, consequently generalizing better OOD.

4 DoCoGen: Domain Counterfactual Generation

We propose a corrupt-and-reconstruct approach for generating D-CONS from given source domain examples (Figure 1). We next extend on these two steps, and describe our filtering mechanism used to disqualify low quality D-CONS.

4.1 Domain Corruption

The first step of generating a D-CON is to mask domain specific terms. In order to mask an example $x \sim \mathcal{D}$ with a destination domain \mathcal{D}' , we first mask all uni-grams w with $m(w, \mathcal{D}, \mathcal{D}') > \tau$, where τ is a hyperparameter and m is a masking score that is defined later in this section. Then, we mask all the remaining bi-grams (that do not contain a masked uni-gram) according to the same masking threshold τ . This process is repeated up to tri-gram expressions. The final output of the corruption step is a masked example $M(x)$.

In Figure 1, the masking scores of uni-grams and bi-grams appear above the input words. An n-gram is masked if and only if its score is above a $\tau = 0.08$ threshold and the scores of its grams are lower. For example, *system* is not masked although the bi-gram *entertainment system* has a score above the τ threshold, since *entertainment* is masked and the score of *system* is lower than τ .

Masking Score Let w be an n-gram and \mathcal{D} be a domain with $n_{\mathcal{D}}$ unlabeled examples. We denote the number of examples from \mathcal{D} that contain w by $\#_{w|\mathcal{D}}$. By assuming that domains have equal prior probabilities and by using the Bayes' rule, the probability of \mathcal{D} given w can be estimated by $P(\mathcal{D} = \mathcal{D} | W = w) \propto \frac{\#_{w|\mathcal{D}} + \alpha}{n_{\mathcal{D}}}$, where α is a smoothing hyperparameter. We define the affinity of w to \mathcal{D} to be:

$$\rho(w, \mathcal{D}) = P(\mathcal{D}|w) \cdot \left(1 - \frac{H(\mathcal{D}|w)}{\log N}\right)$$

where N is the number of unlabeled domains and $H(\mathcal{D}|w)$ is the entropy of $\mathcal{D}|w$, which is upper bounded by $\log N$. Notice that higher $H(\mathcal{D}|w)$ values indicate that w is not related to any specific domain. Finally, we set the masking score of an n-gram w with an origin domain \mathcal{D} and a destination domain \mathcal{D}' as follows:

$$m(w, \mathcal{D}, \mathcal{D}') = \rho(w, \mathcal{D}) - \rho(w, \mathcal{D}')$$

Note that $m(w, \mathcal{D}, \mathcal{D}') \in [-1, 1]$. It can be negative due to the right hand side's subtrahend, which aims to prevent masking n-grams that are related to the destination domain and should appear in the counterfactual, like *system* in Figure 1.

4.2 Domain-Oriented Reconstruction

The second step of DoCoGen is a reconstruction step that involves a generative model, based on

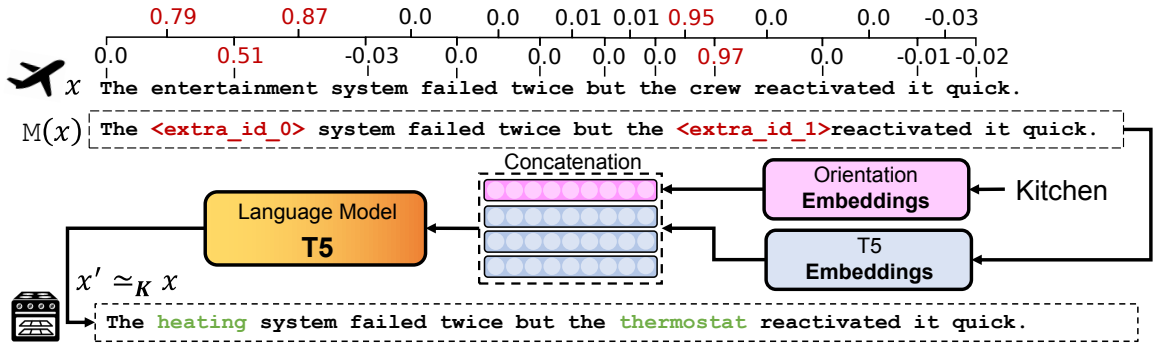


Figure 1: The DoCoGen model. Given a review x from the *airline* domain, we aim to generate a D-CON from the *kitchen* domain. We first corrupt the domain of the example by masking domain specific terms. The numbers above the input words are the masking scores of uni-grams and bi-grams. Terms with scores above a threshold ($\tau = 0.08$) are masked. In the reconstruction step we use a T5-based generation model to generate the D-CON $x' \simeq_K x$. The input of the model is a concatenation of the *orientation vector* that represents the target domain with the model’s embedding vectors which correspond to the tokens of the masked example $M(x)$.

an encoder-decoder T5 architecture (Raffel et al., 2020). Given a masked example $M(x)$ and a destination domain \mathcal{D}' , we concatenate a domain orientation vector v' that represents \mathcal{D}' with the masked input’s embedding vectors. Then, the concatenated matrix is passed as an input to the encoder-decoder model for counterfactual generation, yielding x' . We next describe the mechanism behind domain orientation vectors.

Domain Orientation Vectors In addition to the T5 embedding matrix (T5 Embeddings in Figure 1), we equip our model with another learnable embedding matrix, containing $K \cdot N$ orientation vectors, such that each domain is represented by K different vectors (Orientation Embeddings in Figure 1). We initialize the orientation vectors with the T5 embedding vectors of the domain names and the top $K - 1$ representing words of each domain. The top representing words of domain \mathcal{D} are those which reach the highest score of: $\log(\#_{w|\mathcal{D}} + 1)\rho(w, \mathcal{D})$. We use K orientation vectors to allow us generate a heterogeneous set of D-CONS for a given destination domain (see examples in §A). We note that although the orientation vectors are initialized with vectors from the T5 embedding matrix, they have a different role and thus are likely to converge to different values during the training process.

Training In the spirit of low resource learning, we would like to train DoCoGen in an unsupervised manner, i.e., without access to manually generated D-CONS. Therefore, we use the unlabeled data of our unlabeled domains. For each example x , we provide the model with $M(x)$, the corrupted version of x , and v , the orientation vector of \mathcal{D} , and

with x as the gold output. The model hence learns to reconstruct x given $M(x)$ and v .

Notice that the origin and the destination domains are the same, i.e, $\mathcal{D} = \mathcal{D}'$, and the masking score is $m(w, \mathcal{D}, \mathcal{D}) = 0$. Hence, for masking purposes, we randomly choose $\tilde{\mathcal{D}} \neq \mathcal{D}$ and plug it as the destination domain in the masking score. We then choose an orientation v for \mathcal{D} , by randomly sampling either the domain name or one of its representing words as long as it appears in x .

Finally, since the orientation vector parameters are trained as part of the reconstruction objective, we establish the connection between the orientation vector and the semantics of the completed example. Hence, we expect that at inference time examples will be properly transformed into their D-CONS.

Inference Given $(x, \mathcal{D}, \mathcal{D}')$, we first mask the example to get $M(x)$ and select one orientation vector v' that represents \mathcal{D}' .² Together, the tuple $(M(x), v')$ forms the input, and accordingly the model generates a D-CON $x' \simeq_{\mathcal{D}'} x$. To increase the likelihood that x' originates from \mathcal{D}' , we restrict the model to generate only tokens of the original example or tokens that are related to \mathcal{D}' and meet the condition: $\max_{i \in \{1, \dots, N\}} m(w, \mathcal{D}', \mathcal{D}_i) > \tau$.

4.3 Filtering Mechanism

In order to properly apply DoCoGen within a DA pipeline, we introduce a filtering mechanism that disqualifies low quality D-CONS generated by DoCoGen. Particularly, we train a classifier to predict the domain of the original, human-written unlabeled examples, and use it to remove D-CONS

²§C.2 presents the % of masked tokens in our experiments.

if their predicted domain is not the given destination domain. In addition, we disqualify D-CONS with less than four words or when the word overlap with the original example is lower than 25%. We name DoCoGen when equipped with this filtering mechanism $F\text{-DoCoGen}$.

5 Intrinsic Evaluation

We next assess DoCoGen in terms of its generated D-CONS, ensuring they: (i) belong to the correct domain and label (1, 2), and (ii) are fluent (3, 4). To this end, we collected 20 original reviews, equally distributed among four domains (the A, D, E, and K domains, see §6). We then applied DoCoGen to generate 60 D-CONS, 3 for each of the original reviews (see §6 for the DoCoGen training setup). Finally, we trained the VAE model of Russo et al. (2020) on labeled data (all the labeled data of the A, D, E, and K domains) and applied it to generate five reviews from each of the above four domains, with the same number of positive and negative reviews as in the set of original reviews.

We then conducted a crowd-sourcing experiment where five nearly native English speakers rated each example, considering the following evaluation measures: (1) Domain relevance (D.REL) - whether the topic of the generated text is related to its destination domain; (2) Label preservation (L.PRES) - what is the label of the generated example (and we report whether the answer was identical to the desired label); (3) Linguistic Acceptability (ACCPT) - how logical and grammatical the example is (on a 1-5 scale); and (4) Word error rate (WER) - what is the minimum number of word substitutions, deletions, and insertions that have to be performed to make the example logical and grammatical.³

Table 2 reports our results. DoCoGen achieves high ACCPT scores and low WER scores, significantly outperforming its VAE alternative, which is known to struggle with longer texts (Shen et al., 2019; Iqbal and Qureshi, 2020). Interestingly, DoCoGen achieves compatible results to the original reviews, indicating the high quality of its generated texts. Finally, in more than 90% of the cases DoCoGen manages to change the example domain to the desired domain, and in 80% it preserves the original example label. In comparison, only 88% of the original examples were annotated as their gold label.

³We actually asked the annotators to edit the example and then measured the number of edit operations.

	↑D.REL	↑L.PRES	↑ACCPT	↓WER
VAE	90.0	46.0	2.11	0.54
DoCoGen	93.0	80.0	4.01	0.17
Original Reviews	99.0	88.0	4.73	0.10

Table 2: Human intrinsic evaluation. Up arrows (↑) represent metrics where higher scores are better, and down arrows (↓) represent the opposite.

6 Experimental Setup

6.1 Tasks and Domains⁴

We follow a large body of prior DA work, focusing on the task of binary sentiment classification. Specifically, our experiments include six different domains: the four legacy product review domains (Blitzer et al., 2007) - Books (B), DVDs (D), Electronic items (E) and Kitchen appliances (K); the challenging airline review dataset (A) (Nguyen, 2015; Ziser and Reichart, 2018b); and the restaurant (R) domain obtained from the Yelp dataset challenge (Zhang et al., 2015). The focus of this work is on low resource DA, and thus we randomly sample 100 labeled examples to form the training set for the following domains: A, D, E, and K. Following Ziser and Reichart (2018b), we use 2000 examples for test from each of the target domains and use the following number of unlabeled reviews: A: 39396, D: 34741, E: 13153, and K: 16785.

As described in §2, we explore two DA setups, UDA and ADA. For UDA, where the model has access to unlabeled target domain data, we experiment with a total of 12 cross-domain setups, including the following domains: A, D, E, and K. For ADA, where unlabeled data from the target domain is not within reach, we experiment with a total of 8 setups, including B and R as target domains, and A, D, E, and K as source domains. Our reported results are averaged across 25 different seeds and randomly sampled training and development sets.

DA by Augmentation The DA pipeline includes a T5-based sentiment classifier trained on labeled data from a single source domain and an augmentation model (e.g., DoCoGen) trained on unlabeled data from four unlabeled domains. We first train DoCoGen on the unlabeled data, and then use it for generating D-CONS that enrich the classifier’s training data. For each labeled training example, DoCoGen generates $K = 4$ D-CONS w.r.t. each unlabeled domain, resulting in a total of 16 D-CONS per example. After training the sentiment

⁴URLs of the datasets and the code, implementation and hyperparameter details are described in §B.

	A → D	A → E	A → K	D → A	D → E	D → K	E → A	E → D	E → K	K → A	K → D	K → E	AVG
NoDA	69.4	78.6	78.2	72.3	80.2	82.4	81.0	79.8	87.6	72.5	78.6	85.4	78.8
DANN	70.3	78.7	78.9	75.5	81.2	82.3	82.3	78.3	86.7	81.0	78.3	85.0	79.9
EDA	69.3	79.1	79.4	71.1	79.9	83.0	79.9	80.8	88.0	75.7	80.9	86.4	79.5
RM-RR	69.5	80.1	80.0	72.3	81.0	83.8	79.6	79.5	88.4	70.6	79.1	84.5	79.0
No-OV	67.2	76.5	76.1	71.5	79.7	82.9	80.9	80.5	88.9	74.8	79.6	85.3	78.7
RM-OV	69.3	80.2	80.4	72.7	81.8	84.5	79.6	81.7	89.0	70.3	79.4	85.4	79.5
DoCoGen	70.6	79.7	79.8	75.8	82.8	84.4	83.0	82.0	89.3	81.2	82.2	87.3	81.5
F-DoCoGen	71.1	79.6	79.6	76.7	83.2	84.8	82.6	82.1	89.2	81.4	83.3	88.0	81.8
PERL	72.9	81.1	<u>83.6</u>	81.5	83.0	<u>86.9</u>	81.1	<u>81.7</u>	<u>88.5</u>	77.9	78.2	86.1	81.9
DoCoGen-PERL	<u>75.7</u>	<u>82.7</u>	83.1	<u>82.4</u>	<u>85.0</u>	84.9	<u>81.3</u>	80.8	88.3	<u>79.5</u>	<u>80.9</u>	<u>86.2</u>	<u>82.6</u>
Oracle-Gen	83.8	88.4	88.9	83.6	89.3	90.0	84.9	84.6	90.7	84.1	82.2	89.0	86.6

Table 3: Accuracy scores for each source and target domain pair in the UDA setup. **Bold** numbers mark the best performing T5-based model, and underline numbers mark the best performing PERL-based model.

Source	A		D		E		K		AVG
	B	R	B	R	B	R	B	R	
NoDA	69.1	76.5	82.3	82.8	81.5	84.5	82.4	85.2	80.5
DANN	70.5	77.2	82.7	81.5	80.9	83.4	81.8	83.4	80.2
EDA	69.3	78.0	83.7	82.6	83.2	85.4	82.8	86.3	81.4
RM-RR	69.4	78.4	83.8	83.5	81.9	85.6	83.7	85.4	81.5
No-OV	67.1	76.1	83.8	82.5	82.9	86.2	83.0	85.6	80.9
RM-OV	69.6	78.7	84.3	83.6	83.6	86.2	83.9	85.5	81.9
DoCoGen	70.9	78.1	84.4	82.9	83.9	86.0	84.5	85.7	82.1
F-DoCoGen	71.4	79.3	84.9	83.6	84.2	86.1	85.6	87.2	82.8
Oracle-Gen	84.4	85.2	86.7	86.1	86.0	86.5	85.3	86.5	85.8

Table 4: Accuracy scores for each source and target domain pair in the ADA setup.

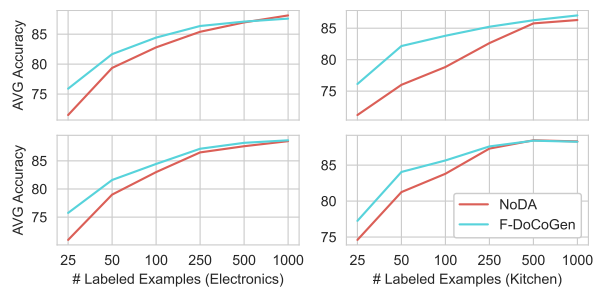


Figure 2: Average accuracy in UDA (top) and ADA (bottom) setups with different number of labeled examples from two source domains: E and K.

classifier on the enriched data, we evaluate it on test examples originating from one of the unlabeled domains (UDA) or one of the unseen domains (ADA). We denote each DA model by the algorithm that was used for enriching its training data.

6.2 Models and Baselines

Our main models are DoCoGen and F-DoCoGen, which is equipped with the filtering mechanism. We compare them to three types of models: (a) baseline models, including both baselines for the entire DA pipeline (1,2,5) and alternative augmentation methods (3,4); (b) ablation models (6,7) that use variants of our D-CON generation algorithm where one component is modified, highlighting the importance of our design choices; and (c) an upper-bound generation model that has access to labeled data from the target domains. Unless otherwise stated, all sentiment classifiers use the same architecture, based on a pre-trained T5 model. We next describe the models in each of these groups.

Baseline DA Models We experiment with five baselines: (1) *No-Domain-Adaptation* (NoDA), A model that is only trained on the available training data from the source domain in each DA setup; (2) *Domain-Adversarial-Neural-Network* (DANN), A model that integrates the sentiment analysis predictive task with an adversarial domain classifier

to learn domain invariant representations (Ganin et al., 2016). This model does not apply augmentation, but instead the unlabeled data is used for training its adversarial component; (3) *Easy-Data-Augmentation* (EDA), an augmentation method that randomly inserts, swaps, and deletes words or replaces synonyms (Wei and Zou, 2019); (4) *Random-masking Random-Reconstructing* (RM-RR), another basic augmentation method that randomly masks tokens from the input example and then fills the masks with tokens that are chosen by a masked language modeling head, as suggested by (Ng et al., 2020); and (5) *PERL*, a SOTA model for the UDA setup (Ben-David et al., 2020).

Ablation Models We consider two variants of DoCoGen: (6) *No-Orientation-Vectors* (No-OV), a generation model that masks tokens by employing a similar masking mechanism as DoCoGen, and then employing a masked language modeling head to fill the masked tokens (without domain orientation vectors); and (7) *Random-Masking with Orientation-Vectors* (RM-OV), a generation model that randomly masks tokens from the input example and then employs the DoCoGen’s reconstruction mechanism to fill the masks.

Upper-Bound We implement an upper-bound model for D-CON augmentation, *Oracle-Matching* (*Oracle-Gen*). Unlike all other models in this work, *Oracle-Gen* has access to target domain labeled data. Thus, given an example from a source domain, *Oracle-Gen* looks for the most similar example with the same label in the target domain, and adds it to its training data (see §B.1).

7 Results

Tables 3 and 4 present the accuracy results for 12 UDA and 8 ADA setups, respectively.

D-CON Generation Impact Our main model, *F-DoCoGen*, outperforms all baseline models (*NoDA*, *DANN*, *EDA*, and *RM-RR*) in 10 of 12 UDA setups and in all ADA setups, exhibiting average performance gains of 1.9% and 1.3% over the best performing baseline model in the UDA (*DANN*) and the ADA (*RM-RR*) setups, respectively. Moreover, *DoCoGen* without filtering, is also superior to all baselines, reaching average gains of 1.6% and of 0.6% across all UDA and ADA setups, respectively. These results highlight the impact of D-CON generation on model robustness to low-resource setups. Finally, our models are also stable: Their std is lower than all baselines (see §C.1).

Ablation Models The tables further demonstrate that *F-DoCoGen* outperforms its ablation models (§ 6.2), namely *No-OV* and *RM-OV*, in 10 of 12 and 7 of 8 UDA and ADA setups, respectively. Furthermore, *F-DoCoGen* achieves an average error reduction of 11.2% and 5.0% in UDA and ADA respectively, over the strongest ablation model (*RM-OV*). Finally, our results demonstrate the importance of inappropriate D-CONS disqualification, as *F-DoCoGen* outperforms *DoCoGen* in 8 of 12 UDA setups and in all ADA setups. This stresses the importance of each of *DoCoGen*’s algorithmic components, i.e. *domain-corruption* (§ 4.1, *F-DoCoGen* vs *RM-OV*), *oriented-reconstruction* (§ 4.2, *F-DoCoGen* vs *No-OV*), and *filtering* (§ 4.3, *F-DoCoGen* vs *DoCoGen*).

Complementary Effect with SOTA Models

We notice that *F-DoCoGen* replicates the average performance of *PERL* (Ben-David et al., 2020), the UDA SOTA. However, since *PERL* is based on a different architecture than the rest of the models (*BERT* vs *T5*), the models are not directly comparable. *PERL* is a pivot-based representation learning

method for DA, which applies pre-training on unlabeled target data and is hence relevant only for UDA. Since *F-DoCoGen* implements a different approach to DA (D-CON generation), we check for the complementary effect of these models: *DoCoGen-PERL* first augments the labeled data with D-CONS and then continues with the *PERL* pipeline. As reported in Table 3, *DoCoGen-PERL* outperforms *PERL* in 8 of 12 UDA setups, providing an average improvement of 0.7%. Furthermore, the average std of *DoCoGen-PERL* is 2.1 compared to 3.6 of *PERL* (§C.1). This stresses the stability of *DoCoGen-PERL* across these challenging setup (Ziser and Reichart, 2019).

Unfortunately, we cannot perform an equivalent comparison in the ADA setup, since its SOTA models (Ben-David et al., 2021; Wright and Augenstein, 2020) employ labeled data from multiple sources. To the best of our knowledge, we are the first to effectively perform single-source ADA.

Training Size Effect We would next like to understand the effect of D-CONS generated by *DoCoGen* on classifiers trained with manually labeled training sets of various sizes. Figure 2 shows that the effect of D-CON augmentation vanishes when the unaugmented classifier reaches accuracy above 85% and a performance plateau (visualized as an elbow in the curve). These results support our hypotheses that low-resource DA scenarios may result in a model that latch on spurious domain correlations, impeding its performance. Accordingly, generating D-CONS by intervening on the domain essentially reduces the reliance on domain-specific information and spurious correlations.

8 Conclusions

We presented *DoCoGen*, a corrupt-and-reconstruct approach for generating domain-counterfactuals (D-CONS) and apply it as a data augmentation method in low-resource DA. We hypothesized that D-CONS may mitigate the reliance on domain-specific features and on spurious correlations and help generalize out of domain.

Our augmentation strategy yields robust models that outperform strong baselines across 20 low-resource sentiment classification DA setups. In future work we would like to further improve the controllable generation quality of *DoCoGen*, potentially extending it to control for multiple attributes. Moreover, we would like our methodology to address additional NLP tasks and DA setups.

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
680
681
682
683
684
685
686
687
688

689
690
691
692

693
694
695
696
697
698

699
700
701

References

Eyal Ben-David, Nadav Oved, and Roi Reichart. 2021. [PADA: A prompt-based autoregressive approach for adaptation to unseen domains](#). *CoRR*, abs/2102.12206.

Eyal Ben-David, Carmel Rabinovitz, and Roi Reichart. 2020. [PERL: pivot-based domain adaptation for pre-trained deep contextualized embedding models](#). *Trans. Assoc. Comput. Linguistics*, 8:504–521.

John Blitzer, Mark Dredze, and Fernando Pereira. 2007. [Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification](#). In *ACL 2007, Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics, June 23-30, 2007, Prague, Czech Republic*. The Association for Computational Linguistics.

John Blitzer, Ryan T. McDonald, and Fernando Pereira. 2006. [Domain adaptation with structural correspondence learning](#). In *EMNLP 2006, Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing, 22-23 July 2006, Sydney, Australia*, pages 120–128. ACL.

Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew M. Dai, Rafal Józefowicz, and Samy Bengio. 2016. [Generating sentences from a continuous space](#). In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany, August 11-12, 2016*, pages 10–21. ACL.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Jiaao Chen, Derek Tam, Colin Raffel, Mohit Bansal, and Diyi Yang. 2021. [An empirical survey of data augmentation for limited data learning in NLP](#). *CoRR*, abs/2106.07499.

Minmin Chen, Zhixiang Eddie Xu, Kilian Q. Weinberger, and Fei Sha. 2012. [Marginalized denoising autoencoders for domain adaptation](#). In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress.

Hal Daumé III and Daniel Marcu. 2006. [Domain adaptation for statistical classifiers](#). *J. Artif. Intell. Res.*, 26:101–126.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. [Understanding back-translation at scale](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 489–500. Association for Computational Linguistics.

Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E. Roberts, Brandon M. Stewart, Victor Veitch, and Diyi Yang. 2021a. [Causal inference in natural language processing: Estimation, prediction, interpretation and beyond](#). *CoRR*, abs/2109.00725.

Amir Feder, Nadav Oved, Uri Shalit, and Roi Reichart. 2021b. [Causalm: Causal model explanation through counterfactual language models](#). *Computational Linguistics*, 47(2):333–386.

Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Edward H. Hovy. 2021. [A survey of data augmentation approaches for NLP](#). In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pages 968–988. Association for Computational Linguistics.

Steven Y. Feng, Aaron W. Li, and Jesse Hoey. 2019. [Keep calm and switch on! preserving sentiment and fluency in semantic text exchange](#). 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 2019, Hong Kong, China, November 3-7, 2019*, pages 2701–2711. Association for Computational Linguistics.

Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor S. Lempitsky. 2016. [Domain-adversarial training of neural networks](#). *The journal of machine learning research*, 17:59:1–59:35.

Matt Gardner, William Merrill, Jesse Dodge, Matthew E. Peters, Alexis Ross, Sameer Singh, and Noah A. Smith. 2021. [Competency problems: On finding and removing artifacts in language data](#). *CoRR*, abs/2104.08646.

Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard S. Zemel, Wieland Brendel,

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
727
728
729

730
731
732
733
734
735
736
737

738
739
740
741
742
743
744
745
746

747
748
749
750
751
752

753
754
755
756
757

758
759

760	Matthias Bethge, and Felix A. Wichmann. 2020. Shortcut learning in deep neural networks . <i>CoRR</i> , abs/2004.07780.	
761		
762		
763	Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. Domain adaptation for large-scale sentiment classification: A deep learning approach . In <i>Proceedings of the 28th International Conference on Machine Learning, ICML 2011, Bellevue, Washington, USA, June 28 - July 2, 2011</i> , pages 513–520. Omnipress.	
764		
765		
766		
767		
768		
769		
770	Xiaochuang Han and Jacob Eisenstein. 2019. Unsupervised domain adaptation of contextualized embeddings for sequence labeling . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 4237–4247. Association for Computational Linguistics.	
771		
772		
773		
774		
775		
776		
777		
778		
779	Jiayuan Huang, Alexander J. Smola, Arthur Gretton, Karsten M. Borgwardt, and Bernhard Schölkopf. 2006. Correcting sample selection bias by unlabeled data . In <i>Advances in Neural Information Processing Systems 19, Proceedings of the Twentieth Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 4-7, 2006</i> , pages 601–608. MIT Press.	
780		
781		
782		
783		
784		
785		
786		
787	Touseef Iqbal and Shaima Qureshi. 2020. The survey: Text generation models in deep learning . <i>Journal of King Saud University-Computer and Information Sciences</i> .	
788		
789		
790		
791	Nitish Joshi and He He. 2021. An investigation of the (in)effectiveness of counterfactually augmented data . <i>CoRR</i> , abs/2107.00753.	
792		
793		
794	Divyansh Kaushik, Eduard H. Hovy, and Zachary Chase Lipton. 2020. Learning the difference that makes A difference with counterfactually-augmented data . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	
795		
796		
797		
798		
799		
800	Divyansh Kaushik, Amrith Setlur, Eduard H. Hovy, and Zachary Chase Lipton. 2021. Explaining the efficacy of counterfactually augmented data . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	
801		
802		
803		
804		
805		
806	Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. CTRL: A conditional transformer language model for controllable generation . <i>CoRR</i> , abs/1909.05858.	
807		
808		
809		
810	Daniel Khashabi, Tushar Khot, and Ashish Sabharwal. 2020. More bang for your buck: Natural perturbation for robust question answering . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 163–170. Association for Computational Linguistics.	
811		
812		
813		
814		
815		
816		
	Sosuke Kobayashi. 2018. Contextual augmentation: Data augmentation by words with paradigmatic relations . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 2 (Short Papers)</i> , pages 452–457. Association for Computational Linguistics.	817
		818
		819
		820
		821
		822
		823
		824
	Ashutosh Kumar, Satwik Bhattamishra, Manik Bhandari, and Partha P. Talukdar. 2019. Submodular optimization-based diverse paraphrasing and its effectiveness in data augmentation . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 3609–3619. Association for Computational Linguistics.	825
		826
		827
		828
		829
		830
		831
		832
		833
		834
	Entony Lekhtman, Yftah Ziser, and Roi Reichart. 2021. DILBERT: customized pre-training for domain adaptation with category shift, with an application to aspect extraction . <i>CoRR</i> , abs/2109.00571.	835
		836
		837
		838
	Jiwei Li, Michel Galley, Chris Brockett, Georgios P. Spithourakis, Jianfeng Gao, and William B. Dolan. 2016. A persona-based neural conversation model . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers</i> . The Association for Computer Linguistics.	839
		840
		841
		842
		843
		844
		845
		846
	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	847
		848
		849
		850
		851
	Clara Meister, Ryan Cotterell, and Tim Vieira. 2020. If beam search is the answer, what was the question? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 2173–2185. Association for Computational Linguistics.	852
		853
		854
		855
		856
		857
	Nathan Ng, Kyunghyun Cho, and Marzyeh Ghassemi. 2020. SSMBA: self-supervised manifold based data augmentation for improving out-of-domain robustness . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 1268–1283. Association for Computational Linguistics.	858
		859
		860
		861
		862
		863
		864
		865
	Quang Nguyen. 2015. The airline review dataset .	866
	Sinno Jialin Pan, Xiaochuan Ni, Jian-Tao Sun, Qiang Yang, and Zheng Chen. 2010. Cross-domain sentiment classification via spectral feature alignment . In <i>Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26-30, 2010</i> , pages 751–760. ACM.	867
		868
		869
		870
		871
		872
		873

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
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

A Additional Generated Examples

We provide additional D-CONS generated by DoCoGen and examples generated by VAE.

(1) D-CONS generated by DoCoGen, for different source-target domain pairs:

(1.1) Original, DVD:

the best show on t.v. Somehow Denis Leary and Peter Tolan managed to improve on an already amazing show. Rescue me continues to be the best show out there.

DVD → Airline:

the best service on board. Somehow staff and flight crew were very good and Tolan managed to improve on an already amazing flight. iberia for me continues to be the best airline out there.

DVD → Electronics:

the best wireless router out there. Somehow ipaq and ethernet tolan managed to improve on an already amazing router. ethernet to me continues to be the best wireless out there.

DVD → Kitchen:

the best oven on the market. Somehow a little bit of work and chef Tolan managed to improve on an already amazing oven. this for me continues to be the best oven out there.

(1.2) Original, Kitchen:

Harder than maple, huh? The very first cut I made on this cutting board left a huge gouge in the surface. Subsequent cuts also left deep marks. For something that's advertised as being harder than maple, I found this to be very surprising. Maybe it's the way it's put together? All I know is that this cutting board is on pace to be the shortest lived ever in my kitchen.

Kitchen → DVD:

Harder than the book, huh? The very first edit I made on this film left a huge gouge in the gut. Subsequent cuts also left deep marks. For something that's advertised as being harder than the book, I found this to be very surprising. Maybe it's the way it's put together? All I know is that this is on pace to be the shortest lived ever in my life.

Kitchen → Electronics:

Harder than a keyboard, huh? The very first click I made on this key board left a huge scratch in the keyboard. Subsequent clicks also left scratch marks. For something that's advertised as being harder than a keyboard, I found this to be very

surprising. Maybe it's the way it's put together? All I know is that this keyboard is on pace to be the shortest lived ever in my life.

(1.3) Original, Airline:

Took flight on a one way basis as return dates were subject to change. I am a top tier AA flier and the addition to Air Berlin gives some good connection options. Aircraft was a new Airbus with standard seat pitch and seat comfort in economy was uncomfortable. The IFE is good with each seat having independent controls and a reasonable selection of movies etc. Unfortunately the selection is via a touch screen buried in your seat.

Airline → DVD:

Took this movie on a trial basis as release dates are subject to change. I am a top ten fan and the addition to the bonus disc gives some good insight. The storyline with standard actors and dialogue in the main role was well done. The plot is well done with each actor having independent controls and a reasonable script. Unfortunately the movie is not about a man who is buried in your soul.

Airline → Electronics:

Took this item on a trial basis as release dates are subject to change. I am a top ten speaker and the addition to the base gives some great options. I have an ipod with standard earbuds and the built in microphone was very nice. The sound is clear with each speaker having independent controls and a reasonable bass response. Unfortunately the remote is not a good buy buried in your pocket.

Airline → Kitchen:

Took this item on a regular basis as delivery dates are subject to change. I am a top cook and the addition to the nonstick gives some good heat retention. I use it with standard pans and the clean up in the dishwasher was easy. The pan is very nice and each pan having independent controls and a reasonable size. Unfortunately the lid is too small for a stovetop and buried in your cupboard.

(2) D-CONS generated by DoCoGen and their masked versions:

(2.1) Original, Airline:

One major problem was that the day I flew there was a fuel service strike so I was delayed for over 4 hours. I frequently check my flight status especially when I leave, throughout that day I had checked at least 4 times. I only found out about it when I

1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178

1179	was at the check in counter. There was no email,	bought other Logitech mouse.	1230
1180	no automated phone call, nothing to notify me. I	<u>Masked text, Electronics → DVD:</u>	1231
1181	was stuck waiting for hours doing nothing.	disappointed with the [MASK] Though I like the	1232
1182	<u>Masked text, Airline → Kitchen:</u>	[MASK], [MASK] two serious problems with	1233
1183	One major problem was that the day I [MASK] a	the [MASK]. First, it is very [MASK] to move-	1234
1184	[MASK] strike so I was [MASK] for [MASK]. I	ment. [MASK] making some [MASK] but it is still	1235
1185	frequently [MASK] my [MASK] especially when	difficult [MASK]. Finally and more importantly,	1236
1186	I leave, [MASK] I had [MASK] at least 4 times. I	the [MASK] about every 8 days. I [MASK] the	1237
1187	[MASK] out about it when I [MASK]. [MASK], no	[MASK] about 6 [MASK] week so I should be get-	1238
1188	[MASK] call, [MASK] me. [MASK] for [MASK]	ting a lot more [MASK] life. I've [MASK] other	1239
1189	doing nothing.	[MASK]	1240
1190	<u>Airline → Kitchen:</u>	<u>Electronics → DVD:</u>	1241
1191	One major problem was that the day I got home	disappointed with the workout. Though I like the	1242
1192	there was a labor strike so I was left without a	workout, I have two serious problems with the	1243
1193	coffee maker for the night. I frequently refill my	workout. First, it is very slow to movement. I'm	1244
1194	coffee especially when I leave, and I had to replace	making some progress but it is still difficult to	1245
1195	it at least 4 times. I found out about it when I got	follow. Finally and more importantly, the workout	1246
1196	home. I sent it back to the store, no one came by	only goes on about every 8 days. I do the workout	1247
1197	to call, and they didn't help me. I sat for over an	about 6 days a week so I should be getting a lot	1248
1198	week doing nothing.	more exercise in my life. I've seen other workouts	1249
1199		that aren't slow.	1250
1200	<u>(2.2) Original, DVD:</u>		1251
1201	The Usual Suspects opened the sealed DVD case	<u>(2.4) Original, Kitchen:</u>	1252
1202	and Get Shorty was inside and not the Usual Sus-	nice cake plate I recieved this cake plate a couple	1253
1203	spects. I sent it back requesting a replacement and	weeks ago. it's very heavy and well made. it came	1254
1204	did not get a replacement. I got a credit, I think.	boxed extra well. The box was inside another box	1255
1205	Not sure. I would rather just have the correct movie.	surrounded by air bags. There was no way this	1256
1206		thing was going to be broken. I am happy with my	1257
1207	<u>Masked text, DVD → Electronics:</u>	purchase.	1258
1208	The [MASK] opened the sealed [MASK] case and	<u>Masked text, Kitchen → Electronics:</u>	1259
1209	Get [MASK] was inside and not the [MASK]. I	nice [MASK] I [MASK] a couple weeks ago. it's	1260
1210	sent it back requesting a replacement and did not	very [MASK]. it came boxed extra well. The box	1261
1211	get a replacement. I got a credit, I think. Not sure.	was inside another box surrounded by air bags.	1262
1212	I would rather just have the correct [MASK].	There was no way this thing was going to be broken.	1263
1213	<u>DVD → Electronics:</u>	I am happy with my purchase.	1264
1214	The router broke when I opened the sealed case	<u>Kitchen → Electronics:</u>	1265
1215	and Get a message that the cable was inside	nice product I bought a couple weeks ago. it's very	1266
1216	and not the router. I sent it back requesting a	easy to use. it came boxed extra well. The box was	1267
1217	replacement and did not get a replacement. I got a	inside another box surrounded by air bags. There	1268
1218	credit, I think. Not sure. I would rather just have	was no way this thing was going to be broken. I	1269
1219	the correct cable.	am happy with my purchase.	1270
1220			1271
1221	<u>(2.3) Original, Electronics:</u>	(3) D-CONS generated by DoCoGen while	1272
1222	disappointed with the mouse Though I like the	providing different orientation vectors:	1273
1223	keyboard, I've found two serious problems with	<u>(3.1) Original, Airline:</u>	1274
1224	the mouse. First, it is very sensitive to movement.	It was a fantastic flight crew helpful and smiling.	1275
1225	I've tried making some adjustments but it is still	All announcements very clear understandable and	1276
1226	difficult to control. Finally and more importantly,	most important things for me multi - language	1277
1227	the mouse burns a set of AA batteries about every	Turkish English Russian and Kazakh. The way	1278
1228	8 days. I only use the mouse about 6 hours a week	aircraft was new.	1279
1229	so I should be getting a lot more battery life. I've		

1280 **Airline** → **DVD**, Orientation: “dvd”:
1281 It was a fantastic movie. All the main actors
1282 were very clear understandable and most impor-
1283 tant things for me multi - language Russian and
1284 Kazakh. The way they talk.

1285 **Airline** → **DVD**, Orientation: “character”:
1286 It was a fantastic movie. All the main character’s
1287 speech was very clear understandable and most
1288 important things for me multi - language Russian
1289 and Kazakh. The way it was done was great.

1290 **Airline** → **DVD**, Orientation: “actor”:
1291 It was a fantastic movie. All actors very clear under-
1292 standable and most important things for me multi -
1293 language Russian and Kazakh. The way they act.

1294 **Airline** → **DVD**, Orientation: “plot”:
1295 It was a fantastic movie. All the plots were very
1296 clear understandable and most important things for
1297 me multi - language Russian and Kazakh. The way
1298 it was told was very good.
1299

1300 (4) Examples generated by the VAE baseline, 1301 controlling both for the domain and the label:

1302 **Airline, Positive:**
1303 i paid for a trip to hawaii. i was traveling with
1304 aeroflot. i was pleasantly surprised. i was pleas-
1305 antly surprised. the service was excellent and ser-
1306 vice. i paid for \$50 for business class. i am very
1307 satisfied with this airline.

1308 **Airline, Positive:**
1309 one of the most memorable movie ever made. i
1310 think this movie is a silly comedy, but i was a little
1311 silly "attitude of" the "buddy" . "attitude" attitude
1312 of the robots, but i was a little silly job of the movie.
1313

1314 **Electronics, Negative:**
1315 not worth the money for my ipod nano. i bought
1316 this product for my 3 year old and i am not sure
1317 why i am not sure why i am not sure why i am not
1318 disappointed.

1319 **Kitchen, Positive:**
1320 broken broken after a broken set of my mother and
1321 i needed a gift for my sister. i was skeptical about
1322 how to do it. i was able to use it to my dishwasher
1323 safe and i was delighted with a silverware. i would
1324 recommend it

B Implementation Details 1325

B.1 URLs of Code and Data 1326

- **DoCoGen Repository** - Code and data will 1327
be released upon acceptance. 1328
- **HuggingFace** (Wolf et al., 2020) - code and 1329
pretrained weights for the T5 model and tok- 1330
enizer: <https://huggingface.co/> 1331
- **SentenceTransformers** (Reimers and 1332
Gurevych, 2019) - code and pretrained 1333
weights of a LM. We use this LM to extract 1334
the embeddings of input examples, and then 1335
calculate the cosine similarity between them 1336
to match examples in the Oracle-Gen 1337
model: <https://www.sbert.net/> 1338
- **PERL** (Ben-David et al., 2020) - A SOTA un- 1339
supervised domain adaptation model: [https:](https://github.com/eyalbd2/PERL) 1340
[://github.com/eyalbd2/PERL](https://github.com/eyalbd2/PERL) 1341
- **NLTK** - code for the Snowball stemmer: 1342
<https://www.nltk.org/index.html> 1343
- **EDA** (Wei and Zou, 2019) - [https://github.](https://github.com/jasonwei20/eda_nlp) 1344
[com/jasonwei20/eda_nlp](https://github.com/jasonwei20/eda_nlp) 1345
- **VAE** - based on the controllable 1346
generation model of Russo et al. 1347
(2020): [https://github.com/DS3Lab/](https://github.com/DS3Lab/control-generate-augment) 1348
[control-generate-augment](https://github.com/DS3Lab/control-generate-augment) 1349

B.2 Hyperparameters and Setups 1350

Data Preprocessing We truncate each example 1351
to 96 tokens, using the HuggingFace T5-base tok- 1352
enizer. The hyper-parameter was set to 96 due to 1353
computation reasons and since the median number 1354
of words in the labeled examples was 89. When an 1355
example is longer than 96 tokens, we keep the first 1356
96 tokens. For examples from the Airline domain, 1357
before truncating, we remove the first sentence 1358
since it mostly contains details about the flight (like 1359
“from JPK to LAX”). 1360

DoCoGen Masking: We estimate $P(\mathcal{D}|w)$ for 1361
uni-grams, bi-grams and tri-grams which appear in 1362
the unlabeled data in at least 10 examples. We use 1363
the NLTK Snowball stemmer to stem each word 1364
of the n-grams. The smoothing hyperparameters in 1365
the computation of $P(\mathcal{D}|w)$ are set to be 1, 5 and 7 1366
for uni-grams, bi-grams and tri-grams, respectively. 1367
We use a $\tau = 0.08$ threshold and mask additional 1368
5% of the training examples (in order to add noise 1369

1370 between training epochs). For RM-RR and RM-OV
1371 we randomly mask 15% of the examples (the stan-
1372 dard ratio for MLM).

1373 **Controllable Model:** We use $K = 4$ orien-
1374 tation vectors for each unlabeled domain and
1375 initialize them with the following representing
1376 words: Airline: {airline, flight, seat,
1377 staff}, DVD: {dvd, character, actor,
1378 plot}, Electronics: {electronics, ipod,
1379 router, software}, Kitchen: {kitchen,
1380 dishwasher, pan, oven}.

1381 The controllable model is based on a pretrained
1382 HuggingFace T5-base model. We train it on
1383 the unlabeled data for 20 epochs and pick the
1384 model whose generated examples for an unlabeled
1385 held-out set are of the highest domain-accuracy
1386 (D.REL).⁵ Training is performed with the AdamW
1387 optimizer (Loshchilov and Hutter, 2019) with a
1388 learning rate parameter of 5e-5 and a weight decay
1389 parameter of 1e-5. For RM-RR and RM-OV we pick
1390 the best models based on a MLM loss computed on
1391 a held-out set. In the example generation step we
1392 use a Beam Search decoding method with a beam
1393 size of 4.

1394 **VAE** As described in the main paper, our VAE
1395 implementation is based on Russo et al. (2020).
1396 To adjust the model for the purposes of this re-
1397 search, we control the task label and the domain
1398 label of each generated review. We train the model
1399 on the entire labeled data and unlabeled data that
1400 is available from four domains: A, D, E, and K,
1401 for a total of 8000 labeled reviews and 104075 un-
1402 labeled reviews. We train the VAE for 60 epochs,
1403 concatenating sentences with more than 96 tokens,
1404 and applying a batch size of 32. The rest of the
1405 hyperparameters were set to the values described
1406 in Russo et al. (2020).

1407 **DA Evaluation** Data Augmentation Given a la-
1408 beled example from the source domain, we gener-
1409 ate $K \cdot N = 16$ examples by DoCoGen, where K
1410 is the number of orientation vectors of each domain
1411 and N is the number of unlabeled domains. We use
1412 the generated examples for data augmentation for
1413 the task classifiers. For all augmentation models,
1414 we apply an augmentation ratio identical to the one
1415 used for DoCoGen, yielding augmented training
1416 sets of the same size. For NoDA and DANN we du-
1417 plicate the training set $K \cdot N$ times, thus the number

⁵The domain accuracy is measured by a domain-classifier trained on the unlabeled data and that is based on the T5 encoder architecture.

of training steps of all the classifiers is identical.
For EDA we use the default hyperparameters.

Sentiment Classifiers All classifiers are based on
the T5-encoder architecture equipped with a linear
layer, except from PERL which is based on the
BERT architecture. We train the classifiers for 5
epochs with a batch size of 64 and pick the best
model based on the label accuracy of the validation
set. Training is performed using the AdamW opti-
mizer with learning rate parameters of 5e-5 for the
encoder blocks and of 5e-4 for the linear layer.

For the results reported in Tables 3, 4, 5 and 6 we
employ a training set that consists of 100 examples
and a validation set with 25 examples. To increase
the robustness of the results in our small labeled
training set setup, we train 25 classifiers, each us-
ing a different randomized seed and a randomly
sampled training set. We report the average perfor-
mance of these classifiers on the test set. For the
results reported in Figure 2, the validation set size
is 25% of the training size. We train the classifiers
on 25 different seeds and partitions for training
sizes 25, 50 and 100, and 10 seeds and partitions
for sizes 250, 500 and 1000.

C Ablation Results

C.1 Standard Deviations

Each of the numbers reported in the main result
tables of the main paper is the average of 25 repeti-
tions, across seeds and training sets. We hence also
report here the standard deviations of these results,
which indicate on the stability of the participating
models.

The standard deviations for the UDA and ADA
setups are presented in Tables 5 and 6, respec-
tively. F-DoCoGen outperforms all baseline mod-
els (NoDA, DANN, EDA, and RM-RR) in 11 of 12
UDA setups and in 6 of 8 ADA setups, demon-
strating a lower average standard deviation and
an improvement of 22.0% and 27.5% in the UDA
and the ADA setups, respectively, over the best
performing baseline model. Moreover, DoCoGen
without filtering is also superior to all baselines.
These results highlight the impact of D-CON gener-
ation on model stability in low-resource DA setups.

As noted in the main paper, we also evalu-
ate the complementary effect of DoCoGen and
PERL, a SOTA model for UDA. Tables 5 shows
that DoCoGen-PERL achieves the lowest aver-
age standard deviation, improving PERL by 42%.
DoCoGen-PERL is hence the best performing

	A → D	A → E	A → K	D → A	D → E	D → K	E → A	E → D	E → K	K → A	K → D	K → E	AVG
NoDA	7.8	6.0	6.8	6.7	5.7	5.4	2.6	4.7	3.0	6.8	4.1	2.9	5.2
DANN	5.4	4.9	5.8	5.2	4.5	4.4	3.1	3.4	3.4	2.8	4.4	2.5	4.1
EDA	6.1	5.7	5.8	7.1	6.8	5.4	4.4	4.9	3.5	6.1	4.5	2.9	5.3
RM-RR	6.8	4.9	5.2	5.7	5.1	4.7	3.2	4.3	2.8	5.5	5.1	3.3	4.7
No-OV	8.0	6.8	7.5	6.8	6.1	5.3	3.0	3.1	2.0	5.0	4.8	3.1	5.1
RM-OV	7.6	4.9	5.4	6.7	5.6	4.7	3.8	2.0	2.0	7.4	4.8	3.1	4.8
DoCoGen	5.9	4.7	5.1	5.5	4.0	3.5	1.9	2.5	2.3	2.2	2.9	1.9	3.5
F-DoCoGen	4.9	4.3	4.8	5.2	3.8	3.1	2.0	2.3	1.9	2.1	2.0	1.7	3.2
PERL	8.3	5.4	4.6	<u>2.0</u>	6.3	<u>1.2</u>	2.3	2.1	<u>0.7</u>	4.7	4.1	1.4	3.6
DoCoGen-PERL	<u>2.2</u>	<u>0.9</u>	<u>2.7</u>	3.0	<u>1.6</u>	2.1	<u>1.9</u>	<u>1.0</u>	2.8	<u>4.1</u>	<u>1.7</u>	<u>0.9</u>	<u>2.1</u>
Oracle-Gen	1.6	1.2	1.7	1.8	1.0	1.4	0.8	1.2	1.0	1.4	2.9	0.9	1.4

Table 5: Standard deviations for each source and target domain pair in the UDA setup. **Bold** numbers mark the best performing T5-based model, and underline numbers mark the best performing PERL-based model.

	A → B	A → R	D → B	D → R	E → B	E → R	K → B	K → R	AVG
NoDA	8.0	6.3	3.5	3.7	5.7	4.0	4.1	2.7	4.8
DANN	6.5	6.2	3.3	3.7	3.3	2.2	3.5	4.2	4.1
EDA	5.9	4.9	4.1	5.0	5.2	4.3	5.0	3.5	4.7
RM-RR	7.0	4.8	2.9	3.5	5.2	2.9	3.5	2.4	4.0
No-OV	8.2	6.2	2.8	4.0	3.7	1.6	4.4	3.1	4.2
RM-OV	7.8	4.9	2.9	4.6	2.6	1.9	3.4	3.3	3.9
DoCoGen	7.0	5.7	2.4	3.4	3.2	1.6	2.6	2.4	3.5
F-DoCoGen	6.0	4.0	2.0	3.3	3.0	1.7	1.9	1.3	2.9
Oracle-Gen	2.1	2.3	2.0	1.6	1.6	1.8	2.4	1.4	1.9

Table 6: Standard deviations for each source and target domain pair in the ADA setup.

↗	A	D	E	K
A	15.2	37.9	37.3	38.0
D	25.0	16.5	24.0	23.9
E	27.8	26.7	15.7	19.7
K	30.2	28.0	21.1	15.7

Table 7: Percents of tokens of the original examples that were masked by DoCoGen. The left column indicates the origin domain and the top row indicates the destination domain.

1468 model both in terms of accuracy (see main paper)
1469 and in terms of standard deviation (stability).

1470 C.2 Masking

1471 Table 7 presents the average percentage of masked
1472 tokens in the corruption step of DoCoGen (see
1473 §4.1). Overall, the average percentage of masked
1474 tokens in a single review is 25.2. These statistics
1475 emphasize the large gap between original reviews
1476 and their D-CONS.