Controlling Neural Machine Translation Formality with Synthetic Supervision

Xing Niu,¹ **Marine Carpuat**² ¹Amazon AWS AI, ²University of Maryland xingniu@amazon.com, marine@cs.umd.edu

Abstract

This work aims to produce translations that convey source language content at a formality level that is appropriate for a particular audience. Framing this problem as a neural sequence-to-sequence task ideally requires training triplets consisting of a bilingual sentence pair labeled with target language formality. However, in practice, available training examples are limited to English sentence pairs of different styles, and bilingual parallel sentences of unknown formality. We introduce a novel training scheme for multi-task models that automatically generates synthetic training triplets by inferring the missing element on the fly, thus enabling endto-end training. Comprehensive automatic and human assessments show that our best model outperforms existing models by producing translations that better match desired formality levels while preserving the source meaning.¹

1 Introduction

Producing language in the appropriate style is a requirement for natural language generation, as the style of a text conveys information beyond its literal meaning (Hovy 1987). This also applies to translation: professional translators adapt translations to their audience (Nida and Taber 2003), yet the output style has been overlooked in machine translation. For example, the French sentence "Bonne idée, mais elle ne convient pas ici." could be translated to "Good idea, but it doesn't fit here.", which is informal because it elides the subject, uses contractions and chained clauses. It could also be translated more formally to "This is a helpful idea. However, it is not suitable for this purpose.", which is grammatically complete and uses more formal and precise terms.

We recently addressed this gap by introducing the task of Formality-Sensitive Machine Translation (FSMT), where given a French sentence and a desired formality level, systems are asked to produce an English translation at the specified formality level (Niu, Martindale, and Carpuat 2017). Building FSMT systems is challenging because of the lack of appropriate training data: bilingual parallel corpora do not come with formality annotations, and parallel corpora of a given provenance do not have a uniform style. Previously, we took a multi-task approach based on sequenceto-sequence models that were trained to perform French-English translation and English formality transfer (Rao and Tetreault 2018) jointly (Niu, Rao, and Carpuat 2018). The resulting multi-task model performs zero-shot FSMT as it has never been exposed to training samples annotated with both reference translation and formality labels.

In this work, we hypothesize that exposing multi-task models to training samples that directly match the FSMT task can help generate formal and informal outputs that differ from each other, and where formality rewrites do not introduce translation errors. We introduce Online Style Inference, an approach to simulate direct supervision by predicting the target formality of parallel sentence pairs on the fly at training time, thus enabling end-to-end training. We also present a variant of side constraints (Sennrich, Haddow, and Birch 2016a) that improves formality control given inputs of arbitrary formality level.²

We conduct a comprehensive automatic and human evaluation of the resulting FSMT systems. First, we show that Online Style Inference introduces more differences between formal and informal translations of the same input, using automatic metrics to quantify lexical and positional differences. Second, we conduct a human evaluation which shows that Online Style Inference preserves the meaning of the input and introduces stronger formality differences compared to a strong baseline. Finally, we analyze the diversity of transformations between formal and informal outputs produced by our approach.

2 A Neural FSMT Model

Neural Machine Translation (NMT) models compute the conditional probability P(Y|X) of translating a source sentence, $X = (x_1, \ldots, x_n)$, to a target sentence, $Y = (y_1, \ldots, y_m)$. By contrast, FSMT requires producing the most likely translation at the given formality level ℓ :

$$\hat{\boldsymbol{Y}} = \operatorname*{arg\,max}_{\boldsymbol{Y}_{\ell}} P(\boldsymbol{Y}_{\ell} \,|\, \boldsymbol{X}, \ell; \boldsymbol{\theta}). \tag{1}$$

Copyright © 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹This work was done when the first author was at the University of Maryland.

²Source code: https://github.com/xingniu/multitask-ft-fsmt.

Ideally, the FSMT model should be trained on triplets $(\boldsymbol{X}, \ell, \boldsymbol{Y}_{\ell})_{1...N}$, but in practice, such training data is not easy to acquire. We tackle this problem by training a crosslingual machine translation model (French→English) and a monolingual bidirectional formality transfer model (Formal-English↔Informal-English) jointly (Niu, Rao, and Carpuat 2018). Specifically, the model is trained on the combination of $(\boldsymbol{X}, \boldsymbol{Y})_{1...N_1}$ and $(\boldsymbol{Y}_{\bar{\ell}}, \ell, \boldsymbol{Y}_{\ell})_{1...N_2}$, where $\boldsymbol{Y}_{\bar{\ell}}$ and \boldsymbol{Y}_{ℓ} have opposite formality levels. The joint model is able to perform zero-shot FSMT by optimizing $\mathcal{L}_{MT} + \mathcal{L}_{FT}$, where

$$\mathcal{L}_{MT} = \sum_{(\boldsymbol{X}, \boldsymbol{Y})} \log P(\boldsymbol{Y} | \boldsymbol{X}; \boldsymbol{\theta}), \qquad (2)$$

$$\mathcal{L}_{FT} = \sum_{(\boldsymbol{Y}_{\bar{\ell}}, \ell, \boldsymbol{Y}_{\ell})} \log P(\boldsymbol{Y}_{\ell} \,|\, \boldsymbol{Y}_{\bar{\ell}}, \ell; \boldsymbol{\theta}). \tag{3}$$

Controlling Output Language Formality

FSMT shares the goal of producing output sentences of a given formality with monolingual formality style transfer tasks. In both cases, the source sentence usually carries its own style and the model should be able to override it with the independent style ℓ . Previously, we achieved this using an attentional sequence-to-sequence model with side constraints (Sennrich, Haddow, and Birch 2016a), i.e., attaching a style tag (e.g., <2Formal>) to the beginning of each source example (Niu, Rao, and Carpuat 2018). In this work, similar to Wang et al. (2018) and Lample et al. (2019), we attach style tags to both source and target sequences to better control output formality given inputs of arbitrary style.

Sennrich, Haddow, and Birch (2016a) hypothesize that source-side tags control the target style because the model "learns to pay attention to the side constraints", but it has not been verified empirically. We hypothesize that the source style tag also influences the encoder hidden states, and providing a target-side tag lets the decoder benefit from encoding style more directly. This approach yields a single model which is able to both transfer formality (e.g., from formal to informal, or vice versa) and preserve formality (e.g., producing an informal output given an informal input).

Synthetic Supervision — Online Style Inference

Prior work on multilingual NMT shows that the translation quality on zero-shot tasks often significantly lags behind when supervision is provided (Johnson et al. 2017). We address this problem by simulating the supervision, i.e., generating synthetic training triplets (X, ℓ, Y) by using the FSMT model itself to predict the missing element of the triplet from parallel sentence pairs (X, Y).

We introduce **Online Style Inference (OSI)** to generate synthetic triplets. Given a translation example (X, Y), we view predicting the formality of Y, i.e., ℓ_Y , as unsupervised classification using only the pre-trained FSMT model.

As illustrated in Figure 1, we use FSMT to produce both informal and formal translations of the same input, $Y_{I} = \text{FSMT}(X, \ell_{I})$ and $Y_{F} = \text{FSMT}(X, \ell_{F})$ respectively.³ We



Figure 1: Online Style Inference. Given a translation example (X, Y), FSMT produces both informal and formal translations of X, i.e., $Y_{I} = \text{FSMT}(X, \ell_{I})$ and $Y_{F} = \text{FSMT}(X, \ell_{F})$. Y is labeled as formal since it is closer to Y_{F} than Y_{I} .

hypothesize that the style of the reference translation Y can be predicted based on its distance from these two translations. For example, if Y is formal, it should be closer to $Y_{\rm F}$ than $Y_{\rm I}$. We measure the closeness by cross-entropy difference (Moore and Lewis 2010, CED), i.e., we calculate the difference of their per-token cross-entropy scores, ${\rm CED}(Y_{\rm I}, Y_{\rm F}) = H_Y(Y_{\rm I}) - H_Y(Y_{\rm F})$. The larger it is, the closer Y is to $Y_{\rm F}$.

Given a positive threshold τ , we label $\ell_{\mathbf{Y}} = \langle 2 \text{Informal} \rangle$ if $\text{CED}(\mathbf{Y}_{\text{I}}, \mathbf{Y}_{\text{F}}) < -\tau$, label $\ell_{\mathbf{Y}} = \langle 2 \text{Formal} \rangle$ if $\text{CED}(\mathbf{Y}_{\text{I}}, \mathbf{Y}_{\text{F}}) > \tau$, and label $\ell_{\mathbf{Y}} = \langle 2 \text{Unknown} \rangle$ otherwise. The threshold τ is chosen dynamically for each mini-batch, and it is equal to the mean of absolute token-level CED of all tokens within a mini-batch. Finally, we are able to generate a synthetic training sample, $(\mathbf{X}, \ell_{\mathbf{Y}}, \mathbf{Y})$, on the fly and optimize $\mathcal{L}_{FT} + \mathcal{L}_{OSI}$, where

$$\mathcal{L}_{OSI} = \sum_{(\boldsymbol{X}, \ell_{\boldsymbol{Y}}, \boldsymbol{Y})} \log P(\boldsymbol{Y} \,|\, \boldsymbol{X}, \ell_{\boldsymbol{Y}}; \boldsymbol{\theta}). \tag{4}$$

Instead of inferring the style label ℓ , we could obtain synthetic training triplets by generating target sequences for a desired ℓ and a given input X. We experiment with one such approach, which we call **Online Target Inference** (OTI), and will compare it with OSI empirically for completeness. However, OTI is expected to be less effective as generates complete output sequences and is therefore more likely to introduce noise in synthetic examples. Given the bilingual parallel sentence pair (X, Y) and a randomly selected target formality ℓ from {<2Informal>, <2Formal>}, we could use the FSMT model to produce a formalityconstrained translation $Y_{\ell}^1 = \text{FSMT}(X, \ell)$. We estimate the quality of Y_{ℓ}^1 indirectly using the multi-task nature of the FSMT models. The FSMT model can also manipulate the formality level of the target side Y via monolingual formality transfer to produce $Y_{\ell}^2 = \text{Transfer}(Y, \ell)$. We hypothesize that the predictions made by these two different paths should be consistent.

The quality of Y_{ℓ}^2 is presumably more reliable than Y_{ℓ}^1 , because the embedded transfer model is trained with direct supervision. We empirically get Y_{ℓ}^2 via greedy search on the fly during the training and use it as the label. Finally, we optimize $\mathcal{L}_{MT} + \mathcal{L}_{FT} + \alpha \mathcal{L}_{OTI}$, where

$$\mathcal{L}_{OTI} = \sum_{(\boldsymbol{X}, \ell, \boldsymbol{Y}_{\ell}^2)} \log P(\boldsymbol{Y}_{\ell}^2 \,|\, \boldsymbol{X}, \ell; \boldsymbol{\theta}).$$
(5)

 $^{{}^{3}}Y_{I}$ and Y_{F} are generated with the *teacher forcing* strategy (Williams and Zipser 1989) given the ground-truth Y.

Corpus	# sentences	# EN tokens
Europarl.v7	1,670,324	39,789,959
News-Commentary.v14	276,358	6,386,435
OpenSubtitles2016	16,000,000	171,034,255

Table 1: Statistics of French-English corpora.

3 Experimental Set-Up

We design experiments to evaluate the impact of our approaches to (1) formality control, and (2) synthetic supervision. We first evaluate formality control on an English style transfer task which provides multiple reference translations to reliably evaluate formality transfer with automatic metrics. We then quantify the differences between formal and informal FSMT translation when using synthetic supervision. Finally, we design a manual evaluation to assess whether synthetic supervision improves over multi-task FSMT. All these experiments share the following set-up.

Data We use the GYAFC corpus introduced by Rao and Tetreault (2018) in all tasks. This corpus consists of informal sentences from Yahoo Answers paired with their formal rewrites by humans. The train split consists of 105K informal-formal sentence pairs whereas the dev/test sets consist of roughly 10K/5K pairs for both formality transfer directions, i.e., $I \rightarrow F$ and $F \rightarrow I$.

We train MT systems on the concatenation of large diverse parallel corpora: (1) Europarl.v7 (Koehn 2005), which is extracted from the proceedings of the European Parliament, and tends to be more formal text; (2) News-Commentary.v14 (Bojar et al. 2018); (3) OpenSubtitles2016 (Lison and Tiedemann 2016), which consists of movie and television subtitles, covers a wider spectrum of styles, but overall tends to be informal since it primarily contains conversations. Following our previous work (Niu, Rao, and Carpuat 2018), we use a bilingual semantic similarity detector to select 16M least divergent examples from $\sim 27.5M$ deduplicated sentence pairs in the original set (Vyas, Niu, and Carpuat 2018).

Preprocessing We apply normalization, tokenization, true-casing, joint source-target BPE with 50,000 operations (Sennrich, Haddow, and Birch 2016b) and sentence-filtering (length 50 cutoff) to parallel training data. Table 1 shows itemized translation data statistics after preprocessing.

Implementation Details We build NMT models upon the attentional RNN encoder-decoder architecture (Bahdanau, Cho, and Bengio 2015) implemented in the Sockeye toolkit (Hieber et al. 2017). Our translation model uses a bidirectional encoder with a single LSTM layer of size 512, multilayer perceptron attention with a layer size of 512, and word representations of size 512. We apply layer normalization (Ba, Kiros, and Hinton 2016), add dropout to embeddings and RNNs (Gal and Ghahramani 2016) with probability 0.2, and tie the source and target embeddings as well as the output layer's weight matrix (Press and Wolf 2017). We

train using the Adam optimizer (Kingma and Ba 2015) with a batch size of 64 sentences and we checkpoint the model every 1000 updates. The learning rate for baseline models is initialized to 0.001 and reduced by 30% after 4 checkpoints without improvement of perplexity on the development set. Training stops after 10 checkpoints without improvement.

We build our FSMT models by fine-tuning the Multi-Task model with the dedicated synthetically supervised objectives described in Section 2, inheriting all settings except the learning rate which is re-initialized to 0.0001. The hyperparameter α in Equation 5 is set to 0.05.

4 Formality Control Evaluation

Our goal is to determine a solid approach for formality control before adding synthetic supervision. For simplicity, we conduct this auxiliary evaluation of formality control on four sub-tasks that use monolingual style transfer data.

Tasks Our task aims to test systems' ability to produce a formal or an informal paraphrase for a given English sentence of arbitrary style, as needed in FSMT. It is derived from formality transfer (Rao and Tetreault 2018), where models transfer sentences from informal to formal $(I \rightarrow F)$ or vice versa $(F \rightarrow I)$. These sub-tasks only evaluate a model's ability in learning mappings between informal and formal languages. We additionally evaluate the ability of systems to preserve formality on informal to informal $(I \rightarrow I)$ and formal to formal $(F \rightarrow F)$ sub-tasks. GYAFC provides four reference target-style human rewrites for each source-style sentences in the test set. For formality preservation, the output is compared with the input sentence in the test set.

Models All models are trained on bidirectional data, which is constructed by swapping the informal and formal sentences of the parallel GYAFC corpus and appending the swapped version to the original. The formality of each target sentence represents the desired input style.

We compare our approach, TAG-SRC-TGT, which attaches tags to both input and output sides, against two baselines. We first implement a baseline method which is trained only on the bidirectional data without showing the target formality (denoted as None). The second baseline is TAG-SRC, the standard method that attaches tags to the source. In addition, we conduct an ablation study on the side constraint method using TAG-SRC-BLOCK, which attaches a tag to the source just like TAG-SRC but blocks the visibility of the tag embeddings from the encoder and retains their connections to the decoder via the attention mechanism (Table 2).

Results Our approach, TAG-SRC-TGT, achieves the best performance overall, reaching the best BLEU scores for three of the four sub-tasks. Comparing with methods acknowledging the target formality (i.e., TAG-SRC*), the None baseline gets slightly lower BLEU scores when it learns to flip the formality on $I \rightarrow F$ and $F \rightarrow I$ tasks.⁴ However, it performs much worse (10-20 BLEU points lower) on

⁴As Rao and Tetreault (2018) note, $F \rightarrow I$ models yield lower BLEU than $I \rightarrow F$ models because informal reference rewrites are

Model	$I \rightarrow F$		$F \rightarrow I$		$I \rightarrow I$		$F \rightarrow F$	
None	70.63 ± 0.23		37.00 ± 0.18		54.54 ± 0.44		58.98 ± 0.93	
TAG-SRC	72.16 ± 0.34	Δ	37.67 ± 0.11	Δ	66.87 ± 0.58	Δ	78.78 ± 0.37	Δ
TAG-SRC-BLOCK	72.00 ± 0.05	-0.16	37.38 ± 0.12	-0.29	65.46 ± 0.29	-1.41	76.72 ± 0.39	-2.06
TAG-SRC-TGT	72.29 ± 0.23	+0.13	37.62 ± 0.37	-0.05	67.81 ± 0.41	+0.94	79.34 ± 0.55	+0.56

Table 2: BLEU scores for variants of side constraint in controlling style on all formality transfer and preservation directions. We report mean and standard deviation over five randomly seeded models. Δ BLEU between each model and the widely used TAG-SRC methods show that (1) blocking the visibility of source tags from the encoder (TAG-SRC-BLOCK) limits its formality control ability; (2) using style tags on both source and target sides (TAG-SRC-TGT) helps control formality better when considering the full range of formality change and formality preservation tasks.

 $I \rightarrow I$ and $F \rightarrow F$ tasks confirming that the None baseline is only able to flip formality and not to preserve it. The TAG-SRC approach is able to preserve formality better than the None baseline, but not as well as TAG-SRC-TGT.

TAG-SRC-BLOCK lags behind TAG-SRC, especially for formality preservation tasks (1-2 BLEU points lower). This discrepancy indicates that the attention mechanism only contributes a portion of the control ability. On the other hand, our proposed variant TAG-SRC-TGT performs better than TAG-SRC on 3/4 tasks (i.e., $I \rightarrow F$, $I \rightarrow I$, and $F \rightarrow F$).

Taken together, these observations show that the impact of tags is not limited to the attention model, and their embeddings influence the hidden representations of encoders and decoders positively. The auxiliary evaluation thus confirms that adding style tags to both source and target sequences is a good approach to model monolingual formality transfer, and therefore motivates using it in our FSMT models as well.

5 Quantifying Differences Between Formal and Informal Outputs in FSMT

Having established the effectiveness of our formality control mechanism, we now turn to the FSMT task and test whether synthetic supervision succeeds in introducing more differences between formal and informal outputs, regardless of translation quality. We will consider translation quality in the next section.

Tasks We test FSMT approaches on two French-English translation test sets with diverse formality: WMT new-stest2014⁵ and MSLT conversation test set⁶. While each test set contains text of varying formality, the written language used in news stories is typically more formal than the spoken language used in conversations.

Baseline Models We start with a standard **NMT** model which is trained with non-tagged French-English parallel data. This model achieves 28.63 BLEU on WMT and 47.83 BLEU on MSLT. We provide these BLEU scores for a sanity check on translation quality.⁷ FSMT models could receive

up to two lower BLEU points on WMT and up to four lower BLEU points on MSLT. However, BLEU cannot be used to evaluate FSMT: given a single reference translation of unknown formality, BLEU penalizes both unwanted translation errors and correct formality rewrites. For example, given the reference "we put together the new wardrobe", the good formal output "we assembled the new wardrobe" and the incorrect output "we accumulated the new wardrobe" would get the same BLEU score.

Next, we compare with Multi-Task, which performs zero-shot FSMT by training machine translation (MT) and formality transfer (FT) jointly. We also compare other two FSMT models introduced in our previous work (Niu, Rao, and Carpuat 2018) for completeness. (1) NMT DS-Tag. It performs data selection on MT training examples (X, Y)using CED in a standard way: it pre-trains language models for informal and formal English in the FT training data and calculates $CED(\mathbf{Y}) = H_{informal}(\mathbf{Y}) - H_{formal}(\mathbf{Y})$. We aim at using all parallel data, for fair comparison, we also conduct three-way tagging as introduced in Section 2. An NMT model is then trained with the formality-tagged training pairs. (2) Multi-Task DS-Tag. It is the combination of Multi-Task and NMT DS-Tag and is trained on both tagged MT pairs and FT pairs. This method is similar to Online Style Inference in terms of tagging training examples using CED. However, Multi-Task DS-Tag uses standard offline language models while Online Style Inference can be interpreted as using source-conditioned online language models.

Metrics Since FSMT quality cannot be evaluated automatically, we devise an approach to quantify surface differences between formal and informal outputs of a given system to guide system development. We define the **Le**xical and **Positional Differences** (LEPOD) score for this purpose, and will come back to FSMT evaluation using human judgments in the next section.

We first compute the pairwise Lexical Difference (LED) based on the percentages of tokens that are not found in both outputs. Formally,

$$\text{LED} = \frac{1}{2} \left(\frac{|S_1 \setminus S_2|}{|S_1|} + \frac{|S_2 \setminus S_1|}{|S_2|} \right),$$
(6)

where S_1 and S_2 is a pair of sequences and $S_1 \setminus S_2$ indicates tokens appearing in S_1 but not in S_2 .

We then compute the pairwise Positional Difference (PoD). (1) We segment the sentence pairs into the longest

highly divergent.

⁵http://www.statmt.org/wmt14/test-full.tgz

⁶https://www.microsoft.com/en-us/download/details.aspx?id= 54689

⁷Detailed BLEU scores are available with released code.

	W	МТ	MSLT					
	LED	PoD	LED	PoD				
NMT	0	0	0	0				
FSMT Baselines								
NMT DS-Tag	9.27	6.44	8.18	1.10				
Multi-Task	10.89	7.76	11.97	1.41				
Multi-Task DS-Tag	11.51	8.35	10.29	1.54				
Multi-Task w/ Synthetic Supervision								
Synth. Target	10.97	7.25	12.40	1.63				
Synth. Style	14.53	12.58	14.52	2.19				

Table 3: LEPOD scores (percentages) show that synthetic supervision introduces more changes between formal and informal outputs than baselines. Online Style Inference (OSI) produces the most diverse informal/formal translations.



Figure 2: Comparing S_1 and S_2 with LEPOD: hollow circles represent non-exact matched tokens, yielding a LED score of $(\frac{7}{15} + \frac{4}{12}) \times \frac{1}{2} = 0.4$. Given the alignment illustrated above, the POD score is $\frac{0+3+2+0}{10} = 0.5$.

sequence of phrasal units that are consistent with the word alignments. Word alignments are obtained using the latest METEOR software (Denkowski and Lavie 2014), which supports stem, synonym and paraphrase matches in addition to exact matches. (2) We compute the maximum distortion within each segment. To do these, we first re-index N aligned words and calculate distortions as the position differences (i.e., index₂ - index₁ in Figure 2). Then, we keep a running total of the distortion array (d_1, d_2, \ldots) , and do segmentation $p = (d_i, \ldots, d_j) \in P$ whenever the accumulation is zero (i.e., $\sum p = 0$). Now we can define

$$POD = \frac{1}{N} \sum_{p \in P} \max(abs(p)).$$
(7)

In extreme cases, when the first word in S_1 is reordered to the last position in S_2 , POD score approaches 1. When words are aligned without any reordering, each alignment constitutes a segment and POD equals 0.

Findings Multi-task methods introduce more differences between formal and informal translations than NMT baselines, and synthetic supervision with Online Target Inference obtains the best lexical and positional difference scores overall (Table 3). Specifically, Multi-Task and Multi-Task DS-Tag get similar lexical and positional variability, and both surpass NMT DS-Tag. Online Target Inference has much larger positional discrepancy scores than all other methods, which indicates that it produces more structural diverse sentences. However, larger surface changes are more likely to alter meaning, and the changes are not guaranteed to be formality-oriented. We therefore turn to human judgments to assess whether meaning is preserved, and whether surface differences are indeed formality related.

6 Human Evaluation of FSMT

Evaluating FSMT systems requires evaluating whether their outputs correctly convey the meaning of the source, and whether the differences between their formal and informal outputs are indicative of formality. Neither LePoD nor BLEU can assess these criteria automatically. We therefore conduct human evaluation to investigate whether synthetic supervision improves over our reimplementation of the state-of-the-art approach (Multi-Task).

Methodology Following Rao and Tetreault (2018) and our previous work (Niu, Rao, and Carpuat 2018), we adopt the following evaluation criteria: *meaning preservation* and *formality difference*.⁸ Our evaluation scheme asks annotators to directly compare sentence pairs on these two criteria and obtains win:tie:loss ratios.

- **Meaning Preservation** We ask annotators to compare outputs of two systems against the reference translation, and decide which one better preserves the reference meaning.
- **Formality Difference** We ask annotators to compare outputs of two systems and decide which is more formal.

We randomly sample \sim 150 examples from WMT and MSLT respectively, and obtain judgments for informal and formal translations of each example. We collect these judgments from 30 volunteers who are native or near-native English speakers. Annotators only compare translations of the same (intended) formality generated by different systems. Identical translation pairs are excluded. Each comparison receives five independent judgments, unless the first three judgments are identical.

The inter-rater agreement using Krippendorff's alpha is ~ 0.5 .⁹ It indicates that there is some variation in annotators' assessment of language formality. We therefore aggregate independent judgments using MACE (Hovy et al. 2013), which estimates the competence of annotators.

Findings Overall, the human evaluation shows that synthetic supervision successfully improves desired formality of the output while preserving translation quality, compared to the multi-task baseline, which represents prior state-of-the-art (Niu, Rao, and Carpuat 2018). Figure 3a and 3b show that Online Style Inference generates informal translations that are annotated as more informal

⁸We do not evaluate fluency because both Rao and Tetreault (2018) and Niu, Rao, and Carpuat (2018) show various automatic systems achieve an almost identical fluency level. Annotators also have systematically biased feeling in fluency when comparing formal and informal sentences (Niu, Martindale, and Carpuat 2017; Rao and Tetreault 2018).

⁹In a sentential formality scoring task, Pavlick and Tetreault (2016) also report relatively low inter-annotator agreement with other measurements.



Figure 3: Win/Tie/Loss counts when comparing Online Style Inference to Multi-Task. Informal translations generated by OSI are annotated as more informal than Multi-Task, while formal translations are annotated as more formal. The OSI model also gets more instances that better preserve the meaning.

Model	identical	contr.	filler	quot.	poss.	y/n	Δ length
Multi-Task	2,140 (33%)	915	530	146	46	13	1.30
Online Target Inference	1,868 (29%)	1,370	635	145	41	21	1.58
Online Style Inference	1,385 (21%)	1,347	530	252	86	33	4.57

Table 4: Heuristic analysis of the differences between informal and formal translations. Synthetic supervision introduce more changes. Online Target Inference usually performs simple substitutions while Online Style Inference performs more less-deterministic changes. Online Style Inference generates more complete and longer formal translations.

(win:tie:loss=151:80:52), while formal translations are annotated as more formal (win:tie:loss=153:84:61). For both cases, the win-loss differences are significant with p < 0.001 using the sign test, where ties are evenly distributed to wins and losses as suggested by Demsar (2006). The results confirm that synthetic supervision lets the model better tailor its outputs to the desired formality, and suggest that the differences between formal and informal outputs detected by the LEPOD scores are indeed representative of formality changes. Figure 3c shows that Online Style Inference preserves the meaning of the source better than Multi-Task (win:tie:loss=205:217:155). The win-loss difference for meaning preservation is still significant with p < 0.02, but is less strong than formality difference.

7 Analysis

How do informal and formal translations differ from each other? Manual inspection reveals that most types of changes made by human rewriters (Pavlick and Tetreault 2016; Rao and Tetreault 2018), including use of filler words, completeness of output and various types of paraphrasing, are observed in our system outputs (see examples in Table 5). We quantify such changes further semi-automatically.

We first check how often formal and informal translations are identical. This happens less frequently with synthetic supervision (Table 4) than with the baseline multi-task system: Online Style Inference system introduces changes between formal and informal translations 12% more often in 6,546 test examples compared to the baseline.

Then, we use rules to check how often simple formality change patterns are found in FSMT outputs (Table 4). A sentence can be made more formal by expanding contractions (contr.) and removing unnecessary fillers such as conjunctions (*so/and/but*) and interjections (*well*) at the beginning

of a sentence (filler). Online Target Inference performs these changes more frequently. We also examine the introduction of quotation marks in formal translations (quot.); using possessive *of* instead of possessive *'s* (poss.); and rewrites of informal use of declarative form for yes-no questions (y/n). Online Style Inference output matches these patterns better than other systems.

Next, we conduct a manual analysis to understand the nature of remaining differences between formal and informal translations of Online Style Inference. We observe that ellipsis is frequent in informal outputs, while formal sentences are more complete, using complement subjects, proper articles, conjunctions, relative pronouns, etc. This is reflected in their longer length (Δ length in Table 4 is the average length difference in characters). Lexical or phrasal paraphrases are frequently used to convey formality, substituting familiar terms with more formal variants (e.g., "grandma" vs. "grandmother"). Examining translations with large POD scores shows that Online Style Inference is more likely to reorder adverbs based on formality: e.g., "I told you already" (I) vs. "I already told you" (F).

A few types of human rewrites categorized by Pavlick and Tetreault (2016) and Rao and Tetreault (2018) are not observed here. For example, our models almost always produce words with correct casing and standard spelling for both informal and formal languages. This matches the characteristics of the translation data we used for training.

Finally, we manually inspect system outputs that fail to preserve the source meaning and reveal some limitations of using synthetic supervision. (1) Inaccurate synthetic labels introduce noise. Online Target Inference sometimes generates "I am not sure" as the formal translation, regardless of the source. We hypothesize that this is due to the imperfect synthetic translations generated by the formality trans-

Туре	Informal translation	Formal translation
Filler	And I think his wife has family there.	I think his wife has family there.
Completeness ▼		
Quotation	The gas tax is simply not sustainable, said Lee.	"The gas tax is simply not sustainable," said Lee.
Yes-No	You like shopping?	Do you like shopping?
Subject	Sorry it's my fault.	I'm sorry it's my fault.
Article	Cookies where I work.	The cookies where I work.
Relativizer	Other stores you can't buy.	The other stores where you can't buy.
Paraphrasing ▼		
Contraction	I think he'd like that, but we'll see.	I think he would like that, but we will see.
Possessive	Fay's innovation perpetuated over the years.	The innovation of Fay has perpetuated over the years.
Adverb	I told you already.	I already told you.
Idiom	Hi, how's it going?	Hi, how are you?
Slang	You gotta let him digest.	You have to let him digest.
Word-1	Actually my dad 's some kind of technician	In fact, my father is some kind of technician
	so he understands, but my mom 's very old.	so he understands, but my mother is very old.
Word-2	Maybe a little more in some areas.	Perhaps a little more in certain areas.
Word-3	It's really necessary for our nation.	This is essential for our nation.
Phrase-1	Yeah, me neither .	Yeah, neither do I.
Phrase-2	I think he's moving to California now.	I think he is moving to California at the moment.
Phrase-3	It could be a Midwest thing.	This could be one thing from the Midwest.

Table 5: Range of differences between informal and formal translations from the Online Style Inference model output.

fer sub-model reinforce this error pattern. (2) Synthetic data may not reflect the true distribution. Occasionally, Online Style Inference drops the first word in a formal sentence even if it is not a filler, e.g. "On Thursday, …" We hypothesize that labeling too many formal/informal examples of similar patterns could lead to ignoring context. While Online Style Inference improves meaning preservation comparatively, it still bears the challenge of altering meaning when fitting to a certain formality, such as generating "there will be no longer than the hill of Runyonyi" when the reference is "then only Rumyoni hill will be left".

8 Related Work

Controlling the output style in MT has received sparse attention. The pioneering work by Mima, Furuse, and Iida (1997) improves rule-based MT using extra-linguistic information such as speaker's role and gender. With the success of statistical MT models, people usually define styles by selecting representative data. After pre-selecting relevant data offline, Lewis, Federmann, and Xin (2015) and van der Wees, Bisazza, and Monz (2016) build conversational MT systems, Rabinovich et al. (2017) and Michel and Neubig (2018) build personalized (gender-specific) MT systems, Sennrich, Haddow, and Birch (2016a) control the output preference of T-V pronouns, while Yamagishi et al. (2016) control the active/passive voice of the translation. In contrast, we dynamically generate synthetic supervision and our methods outperform offline data selection.

Multi-task FSMT is closely related to zero-shot multilingual NMT. Johnson et al. (2017) first built a multilingual NMT system using shared NMT encoder-decoders for all languages with target language specifiers. The resulting system can translate between language pairs that are never trained on, but performs worse than supervised models and even the simple pivoting approach for those language pairs. Strategies to mitigate this problem include target word filtering (Ha, Niehues, and Waibel 2017), dedicated attention modules (Blackwood, Ballesteros, and Ward 2018), generating dedicated encoder-decoder parameters (Platanios et al. 2018) and encouraging the model to use source-language invariant representations (Arivazhagan et al. 2019). We address this problem from a different perspective for the FSMT task by automatically inferring style labels. Our Online Target Inference approach is similar in spirit to a contemporaneous method that encourages the model to produce equivalent translations of parallel sentences into an auxiliary language (Al-Shedivat and Parikh 2019).

9 Conclusion

This paper showed that synthetic supervision improves multi-task models for formality-sensitive machine translation. We introduced a novel training scheme for multi-task models that, given bilingual parallel examples and monolingual formality transfer examples, automatically generate synthetic training triples by inferring the target formality from a given translation pair. Human evaluation shows that this approach outperforms a strong multi-task baseline by producing translations that better match desired formality levels while preserving the source meaning. Additional automatic evaluation shows that (1) attaching style tags to both input and output sequences improves the ability of a single model to control formality, by not only transferring but also preserving formality when required; and (2) synthetic supervision via Online Target Inference introduces more changes between formal and informal translations of the same input. Analysis shows that these changes span almost all types of changes made by human rewriters.

Taken together, these results show the promise of syn-

thetic supervision for style-controlled language generation applications. In future work, we will investigate scenarios where style transfer examples are not readily available, including for languages other than English, and for style distinctions that are more implicit and not limited to binary formal-informal distinctions.

References

Al-Shedivat, M., and Parikh, A. 2019. Consistency by agreement in zero-shot neural machine translation. In *Proceedings of NAACL*, 1184–1197.

Arivazhagan, N.; Bapna, A.; Firat, O.; Aharoni, R.; Johnson, M.; and Macherey, W. 2019. The missing ingredient in zero-shot neural machine translation. *CoRR* abs/1903.07091.

Ba, L. J.; Kiros, R.; and Hinton, G. E. 2016. Layer normalization. *CoRR* abs/1607.06450.

Bahdanau, D.; Cho, K.; and Bengio, Y. 2015. Neural machine translation by jointly learning to align and translate. In *Proceedings* of *ICLR*.

Blackwood, G. W.; Ballesteros, M.; and Ward, T. 2018. Multilingual neural machine translation with task-specific attention. In *Proceedings of COLING*, 3112–3122.

Bojar, O.; Federmann, C.; Fishel, M.; Graham, Y.; Haddow, B.; Koehn, P.; and Monz, C. 2018. Findings of the 2018 conference on machine translation (WMT18). In *Proceedings of WMT*, 272–303.

Demsar, J. 2006. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research* 7:1–30.

Denkowski, M. J., and Lavie, A. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of WMT*, 376–380.

Ficler, J., and Goldberg, Y. 2017. Controlling linguistic style aspects in neural language generation. In *Proceedings of the Workshop on Stylistic Variation*, 94–104.

Fu, Z.; Tan, X.; Peng, N.; Zhao, D.; and Yan, R. 2018. Style transfer in text: Exploration and evaluation. In *Proceedings of AAAI*, 663–670.

Gal, Y., and Ghahramani, Z. 2016. A theoretically grounded application of dropout in recurrent neural networks. In *Advances in Neural Information Processing Systems* 29, 1019–1027.

Ha, T.; Niehues, J.; and Waibel, A. H. 2017. Effective strategies in zero-shot neural machine translation. In *Proceedings of IWSLT*.

Hieber, F.; Domhan, T.; Denkowski, M.; Vilar, D.; Sokolov, A.; Clifton, A.; and Post, M. 2017. Sockeye: A toolkit for neural machine translation. *CoRR* abs/1712.05690.

Hovy, D.; Berg-Kirkpatrick, T.; Vaswani, A.; and Hovy, E. H. 2013. Learning whom to trust with MACE. In *Proceedings of NAACL*, 1120–1130.

Hovy, E. 1987. Generating natural language under pragmatic constraints. *Journal of Pragmatics* 11(6):689–719.

Hu, Z.; Yang, Z.; Liang, X.; Salakhutdinov, R.; and Xing, E. P. 2017. Toward controlled generation of text. In Precup, D., and Teh, Y. W., eds., *Proceedings of ICML*, volume 70, 1587–1596.

Johnson, M.; Schuster, M.; Le, Q. V.; Krikun, M.; Wu, Y.; Chen, Z.; Thorat, N.; Viégas, F. B.; Wattenberg, M.; Corrado, G.; Hughes, M.; and Dean, J. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics* 5:339–351.

Kingma, D. P., and Ba, J. 2015. Adam: A method for stochastic optimization. In *Proceedings of ICLR*.

Koehn, P. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of the Tenth Machine Translation Summit*, 79–86.

Korotkova, E.; Del, M.; and Fishel, M. 2018. Monolingual and cross-lingual zero-shot style transfer. *CoRR* abs/1808.00179.

Lample, G.; Subramanian, S.; Smith, E.; Denoyer, L.; Ranzato, M.; and Boureau, Y.-L. 2019. Multiple-attribute text rewriting. In *Proceedings of ICLR*.

Lewis, W.; Federmann, C.; and Xin, Y. 2015. Applying crossentropy difference for selecting parallel training data from publicly available sources for conversational machine translation. In *Proceedings of IWSLT*.

Lison, P., and Tiedemann, J. 2016. Opensubtitles2016: Extracting large parallel corpora from movie and TV subtitles. In *Proceedings of LREC*.

Luong, T.; Pham, H.; and Manning, C. D. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of EMNLP*, 1412–1421.

Michel, P., and Neubig, G. 2018. Extreme adaptation for personalized neural machine translation. In *Proceedings of ACL*, 312–318.

Mima, H.; Furuse, O.; and Iida, H. 1997. Improving performance of transfer-driven machine translation with extra-linguistic informatioon from context, situation and environment. In *Proceedings of IJCAI*, 983–989.

Moore, R. C., and Lewis, W. D. 2010. Intelligent selection of language model training data. In *Proceedings of ACL*, 220–224.

Mueller, J.; Gifford, D. K.; and Jaakkola, T. S. 2017. Sequence to better sequence: Continuous revision of combinatorial structures. In *Proceedings of ICML*, volume 70, 2536–2544.

Nida, E. A., and Taber, C. R. 2003. *The Theory and Practice of Translation*, volume 8 of *Helps for Translators*. Brill.

Niu, X.; Martindale, M. J.; and Carpuat, M. 2017. A study of style in machine translation: Controlling the formality of machine translation output. In *Proceedings of EMNLP*, 2814–2819.

Niu, X.; Rao, S.; and Carpuat, M. 2018. Multi-task neural models for translating between styles within and across languages. In *Proceedings of COLING*, 1008–1021.

Pavlick, E., and Tetreault, J. R. 2016. An empirical analysis of formality in online communication. *Transactions of the Association for Computational Linguistics* 4:61–74.

Platanios, E. A.; Sachan, M.; Neubig, G.; and Mitchell, T. M. 2018. Contextual parameter generation for universal neural machine translation. In *Proceedings of EMNLP*, 425–435.

Press, O., and Wolf, L. 2017. Using the output embedding to improve language models. In *Proceedings of EACL*, 157–163.

Rabinovich, E.; Mirkin, S.; Patel, R. N.; Specia, L.; and Wintner, S. 2017. Personalized machine translation: Preserving original author traits. In *Proceedings of EACL*, 1074–1084.

Rao, S., and Tetreault, J. R. 2018. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In *Proceedings of NAACL*, 129–140.

Sennrich, R., and Haddow, B. 2016. Linguistic input features improve neural machine translation. In *Proceedings of the First Conference on Machine Translation*, 83–91.

Sennrich, R.; Haddow, B.; and Birch, A. 2016a. Controlling politeness in neural machine translation via side constraints. In *Proceedings of NAACL*, 35–40.

Sennrich, R.; Haddow, B.; and Birch, A. 2016b. Neural machine translation of rare words with subword units. In *Proceedings of ACL*, 1715–1725.

Shen, T.; Lei, T.; Barzilay, R.; and Jaakkola, T. S. 2017. Style transfer from non-parallel text by cross-alignment. In *Advances in Neural Information Processing Systems 30*, 6833–6844.

van der Wees, M.; Bisazza, A.; and Monz, C. 2016. Measuring the effect of conversational aspects on machine translation quality. In *Proceedings of COLING*, 2571–2581.

Vyas, Y.; Niu, X.; and Carpuat, M. 2018. Identifying semantic divergences in parallel text without annotations. In *Proceedings of NAACL*, 1503–1515.

Wang, Y.; Zhang, J.; Zhai, F.; Xu, J.; and Zong, C. 2018. Three strategies to improve one-to-many multilingual translation. In *Proceedings of EMNLP*, 2955–2960.

Williams, R. J., and Zipser, D. 1989. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation* 1(2):270–280.

Xu, K.; Ba, J.; Kiros, R.; Cho, K.; Courville, A. C.; Salakhutdinov, R.; Zemel, R. S.; and Bengio, Y. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *Proceedings* of *ICML*, volume 37, 2048–2057.

Yamagishi, H.; Kanouchi, S.; Sato, T.; and Komachi, M. 2016. Controlling the voice of a sentence in japanese-to-english neural machine translation. In *Proceedings of the 3rd Workshop on Asian Translation*, 203–210.

A Details of the Human Evaluation

As described in the main paper, we assess model outputs on two criteria: *meaning preservation* and *formality difference*.

Meaning Preservation The following instruction is provided to annotators.

For each task, you will be presented with an English sentence and two rewrites of that sentence. Your task is to judge which rewrite better preserves the meaning of the original and choose from:

- Rewrite 1 is much better
- Rewrite 1 is better
- No preference between Rewrite 1 and Rewrite 2 (no difference in meaning or hard to say)
- Rewrite 2 is better
- Rewrite 2 is much better

Note that this task focuses on differences in content, so differences in style (such as formality) between the original and rewrites are considered okay. [Some examples with explanations are provided.]

Formality Difference The following instruction is provided to annotators.

People use different varieties of language depending on the situation: formal language is required in news articles, official speeches or academic assignments, while informal language is more appropriate in instant messages or spoken conversations between friends.

You will be presented with two English sentences, and your task is to decide which one is more formal and choose from:

- Sentence 1 is much more formal
- Sentence 1 is more formal
- No preference between Sentence 1 and Sentence 2 (no difference in formality or hard to say)
- Sentence 2 is more formal
- Sentence 2 is much more formal

Keep in mind:

- Language formality can be affected by many factors, such as the choices of grammar, vocabulary, and punctuation.
- The sentences in the pair could have different meanings. Please rate the formality of the sentences independent of their meaning.
- The sentences in the pair could be nonsensical. Please rate the formality of the sentences independent of their quality.

Generally, a sentence with small formality changes such as fewer contractions, proper punctuation or some formal terms is considered "more formal". A sentence is considered "much more formal" if it contains multiple indicators of formality, or if the sentence construction itself reflects a more formal style. That said, feel free to use your own judgment for doing the task if what you see is not covered by these examples. [Some examples with explanations are provided.]

B Extended Formality Control Evaluation

While the implementations of neural language generation converge to an encoder-decoder framework, the design choice of controlling the style is full of variety. The style information could be injected into different parts of the heterogeneous neural network and all roads lead to Rome. However, those implementations have never been compared with a controlled experiment and analyzed contrastively. We therefore conduct a benchmark test on a four-way formality rewriting task introduced in the main paper.

In order to focus on the designing of style-sensitive neural models, we compare methods performing formality rewriting with a single encoder and a single decoder, in contrast with a dedicated model or decoder for each transfer direction (Rao and Tetreault 2018; Fu et al. 2018). The attention mechanism is the *de facto* standard for language generation (Bahdanau, Cho, and Bengio 2015; Luong, Pham, and Manning 2015; Xu et al. 2015), we therefore compare methods being compatible with the attention mechanism, in contract with methods that compress the content into one single vector (Mueller, Gifford, and Jaakkola 2017; Hu et al. 2017; Shen et al. 2017; Fu et al. 2018).

We compare the following methods, with a focus on their complexity and effectiveness.

- **TAG-SRC-TGT** This is the method used in the main paper. It attaches style tags to both source and target sequences. Each additional tag occupies one embedding vector, which has a size of $O(E_w)$, where E_w is the word embedding size.
- **FACTOR** The style information can be incorporated as source word factors, which is implemented as style factor embeddings concatenated to the word embeddings (Sennrich and Haddow 2016), i.e., $\tilde{X}_i = [X_i; X_i^{\text{style}}]$. Korotkova, Del, and Fishel (2018) adopt this design choice for multiple-style transfer due to its flexibility. Summing the factor and word embeddings of the same size is another combination strategy and we name it FACTOR-SUM, which uses $O(E_w)$ space, as opposite to FACTOR-CONCAT, which uses $O(E_s)$ space. E_s is the style embedding size and usually much smaller than E_w .¹⁰
- **PRED** Alternatively, we can inject the style information later on to the decoder by concatenating style embeddings to predicted target word embeddings, i.e., $\tilde{Y}_t = [Y_t; Y_t^{\text{style}}]$. Ficler and Goldberg (2017) use this method in the style-conditioned language

¹⁰We use $E_s = 5$ in our experiments.

Model	$I \rightarrow F$		$F \rightarrow I$		$I \rightarrow I$		$F \rightarrow F$	
TAG-SRC-TGT	72.29 ± 0.23	Δ	37.62 ± 0.37	Δ	67.81 ± 0.41	Δ	79.34 ± 0.55	Δ
FACTOR-CONCAT	72.47 ± 0.11	+0.18	37.62 ± 0.26	0.00	67.03 ± 0.36	-0.78	79.80 ± 0.38	+0.46
FACTOR-SUM	72.43 ± 0.29	+0.14	37.78 ± 0.26	+0.16	67.24 ± 0.56	-0.57	80.34 ± 0.46	+1.00
Pred-Concat	72.35 ± 0.16	+0.06	37.62 ± 0.13	0.00	66.69 ± 0.21	-1.12	77.85 ± 0.31	-1.49
Pred-Sum	72.02 ± 0.30	-0.27	37.41 ± 0.17	-0.21	66.15 ± 0.41	-1.66	77.62 ± 0.28	-1.72
BOS	72.08 ± 0.22	-0.21	37.56 ± 0.13	-0.06	66.40 ± 0.23	-1.41	77.43 ± 0.34	-1.91
BIAS	71.58 ± 0.31	-0.71	37.52 ± 0.15	-0.10	63.66 ± 0.51	-4.15	73.24 ± 0.55	-6.10

Table 6: BLEU scores of various methods for controlling the style on four formality transfer (preservation) directions. The numbers before and after ' \pm ' are the mean and standard deviation over five randomly seeded models. Methods are compared with TAG-SRC-TGT and the Δ BLEU scores are listed.

generation. Summing can also be used here, and we name these two variants PRED-CONCAT and PRED-SUM. The space complexities are $O(E_s)$ and $O(E_w)$.

Note that for both FACTOR and PRED, the computational complexity also increase proportionally with the sequence length since the style embeddings are combined with word embeddings for each time step.

- **BOS** Analogical to TAG, the target style embeddings can be dynamically attached to the target sequence as a begin-of-sequence symbol (<BOS>). This approach has been successfully applied to multiple-attribute text rewriting (Lample et al. 2019). Each stylistic <BOS> embedding occupies $O(E_w)$ space.
- **BIAS** The bias parameter influences the model's lexical choice in the output layer (i.e. $\operatorname{softmax}(\boldsymbol{W}\boldsymbol{h}_t + \boldsymbol{b})$), so we can assign a dedicated bias for each style. Michel and Neubig (2018) use this technique in personalized NMT. Each dedicated bias has a size of O(V), where V is the vocabulary size.

We compare all methods to TAG-SRC-TGT, which is introduced in the main paper. Experimental settings and implementation details are identical to the intrinsic evaluation in the main paper and we report scores in Table 6.

The other method family incorporating the style information as early as at the encoding stage, FACTOR, waxes and wanes, and performs similar to TAG-SRC-TGT. Remaining methods that incorporate the style information only to the decoder, on the other hand, get lower BLEU scores across the board.