Unsupervised Text Summarization via Mixed Model Back-Translation

Yacine Jernite

Facebook AI Research, New York, NY yjernite@fb.com

Abstract

Back-translation based approaches have recently lead to significant progress in unsupervised sequence-to-sequence tasks such as machine translation or style transfer. In this work, we extend the paradigm to the problem of learning a sentence summarization system from unaligned data. We present several initial models which rely on the asymmetrical nature of the task to perform the first back-translation step, and demonstrate the value of combining the data created by these diverse initialization methods. Our system outperforms the current state-of-the-art for unsupervised sentence summarization from fully unaligned data by over 2 ROUGE, and matches the performance of recent semi-supervised approaches.

1 Introduction

Machine summarization systems have made significant progress in recent years, especially in the domain of news text. This has been made possible among other things by the popularization of the neural sequence-to-sequence (seq2seq) paradigm (Kalchbrenner and Blunsom, 2013; Sutskever et al., 2014; Cho et al., 2014), the development of methods which combine the strengths of extractive and abstractive approaches to summarization (See et al., 2017; Gehrmann et al., 2018), and the availability of large training datasets for the task, such as Gigaword or the CNN-Daily Mail corpus which comprise of over 3.8M shorter and 300K longer articles and aligned summaries respectively. Unfortunately, the lack of datasets of similar scale for other text genres remains a limiting factor when attempting to take full advantage of these modeling advances using supervised training algorithms.

In this work, we investigate the application of back-translation to training a summarization system in an unsupervised fashion from unaligned full text and summaries corpora. Back-translation has been successfully applied to unsupervised training for other sequence to sequence tasks such as machine translation (Lample et al., 2018) or style transfer (Subramanian et al., 2018). We outline the main differences between these settings and text summarization, devise initialization strategies which take advantage of the asymmetrical nature of the task, and demonstrate the advantage of combining varied initializers. Our approach outperforms the previous state-of-the-art on unsupervised text summarization while using less training data, and even matches the ROUGE scores of recent semi-supervised methods.

2 Related Work

Rush et al. (2015)'s work on applying neural seq2seq systems to the task of text summarization has been followed by a number of works improving upon the initial model architecture. These have included changing the base encoder structure (Chopra et al., 2016), adding a pointer mechanism to directly re-use input words in the summary (Nallapati et al., 2016; See et al., 2017), or explicitly pre-selecting parts of the full text to focus on (Gehrmann et al., 2018). While there have been comparatively few attempts to train these models with less supervision, auto-encoding based approaches have met some success (Miao and Blunsom, 2016; Wang and Lee, 2018).

Miao and Blunsom (2016)'s work endeavors to use summaries as a discrete latent variable for a text auto-encoder. They train a system on a combination of the classical log-likelihood loss of the supervised setting and a reconstruction objective which requires the full text to be mostly recoverable from the produced summary. While their method is able to take advantage of unlabelled data, it relies on a good initialization of the encoder part of the system which still needs

to be learned on a significant number of aligned pairs. Wang and Lee (2018) expand upon this approach by replacing the need for supervised data with adversarial objectives which encourage the summaries to be structured like natural language, allowing them to train a system in a fully unsupervised setting from unaligned corpora of full text and summary sequences. Finally, (Song et al., 2019) uses a general purpose pre-trained text encoder to learn a summarization system from fewer examples. Their proposed MASS scheme is shown to be more efficient than BERT (Devlin et al., 2018) or Denoising Auto-Encoders (DAE) (Vincent et al., 2008; Fu et al., 2018).

This work proposes a different approach to unsupervised training based on back-translation. The idea of using an initial weak system to create and iteratively refine artificial training data for a supervised algorithm has been successfully applied to semi-supervised (Sennrich et al., 2016) and unsupervised machine translation (Lample et al., 2018) as well as style transfer (Subramanian et al., 2018). We investigate how the same general paradigm may be applied to the task of summarizing text.

3 Mixed Model Back-Translation

Let us consider the task of transforming a sequence in domain A into a corresponding sequence in domain B (e.g. sentences in two languages for machine translation). Let \mathcal{D}_A and \mathcal{D}_B be corpora of sequences in A and B, without any mapping between their respective elements. The back-translation approach starts with initial seq2seq models $f_{A\to B}^0$ and $f_{B\to A}^0$, which can be hand-crafted or learned without aligned pairs, and uses them to create artificial aligned training data:

$$\mathcal{D}_{A'\to B}^0 = \left\{ \left(f_{B\to A}^0(b), b \right); \ \forall b \in \mathcal{D}_B \right\}$$
 (1)

$$\mathcal{D}_{B'\to A}^0 = \left\{ \left(f_{A\to B}^0(a), a \right); \ \forall a \in \mathcal{D}_A \right\} \quad (2)$$

Let S denote a supervised learning algorithm, which takes a set of aligned sequence pairs and returns a mapping function. This artificial data can then be used to train the next iteration of seq2seq models, which in turn are used to create new artificial training sets (A and B can be switched here):

$$f_{A \to B}^{i+1} = \mathcal{S}(\mathcal{D}_{A' \to B}^i) \tag{3}$$

$$\mathcal{D}_{B'\to A}^{i+1} = \left\{ \left(f_{A\to B}^{i+1}(a), a \right); \ \forall a \in \mathcal{D}_A \right\}$$
 (4)

The model is trained at each iteration on artificial inputs and real outputs, then used to create new

training inputs. Thus, if the initial system isn't too far off, we can hope that training pairs get closer to the true data distribution with each step, allowing in turn to train better models.

In the case of summarization, we consider the domains of full text sequences \mathcal{D}^F and of summaries \mathcal{D}^S , and attempt to learn *summarization* $(f_{F \to S})$ and *expansion* $(f_{S \to F})$ functions. However, contrary to the translation case, \mathcal{D}^F and \mathcal{D}^S are not interchangeable. Considering that a summary typically has less information than the corresponding full text, we choose to only define initial $F \to S$ models. We can still follow the proposed procedure by alternating directions at each step.

3.1 Initialization Models for Summarization

To initiate their process for the case of machine translation, Lample et al. (2018) use two different initialization models for their neural (NMT) and phrase-based (PBSMT) systems. The former relies on denoising auto-encoders in both languages with a shared latent space, while the latter uses the PBSMT system of Koehn et al. (2003) with a phrase table obtained through unsupervised vocabulary alignment as in (Grave et al., 2018).

While both of these methods work well for machine translation, they rely on the input and output having similar lengths and information content. In particular, the statistical machine translation algorithm tries to align most input tokens to an output word. In the case of text summarization, however, there is an inherent asymmetry between the full text and the summaries, since the latter express only a subset of the former. Next, we propose three initialization systems which implicitly model this information loss. Full implementation details are provided in the Appendix.

Procrustes Thresholded Alignment (Pr-Thr)

The first initialization is similar to the one for PB-SMT in that it relies on unsupervised vocabulary alignment. Specifically, we train two skipgram word embedding models using FASTTEXT (Bojanowski et al., 2017) on \mathcal{D}^F and \mathcal{D}^S , then align them in a common space using the Wasserstein Procrustes method of Grave et al. (2018). Then, we map each word of a full text sequence to its nearest neighbor in the aligned space if their distance is smaller than some threshold, or skip it otherwise. We also limit the output length, keeping only the first N tokens. We refer to this function as $f_E^{(\text{Pr-Thr}),0}$.

(Original) france took an important step toward power market liberalization monday, braving union anger to announce the partial privatization of state-owned behemoth electricite de france.

(Pr-Thr) france launched a partial UNK of state-controlled utility, the privatization agency said.

(DBAE) france's state-owned gaz de france sa said tuesday it was considering partial partial privatization of france's state-owned nuclear power plants.

 $(\mu:1)$ france launches an initial public announcement wednesday as the european union announced it would soon undertake a partial privatization.

(Title) france launches partial edf privatization

Table 1: Full text sequences generated by $f_{S \to F}^{(\text{Pr-Thr}),1}$, $f_{S \to F}^{(\text{DBAE}),1}$, and $f_{S \to F}^{(\boldsymbol{\mu}:1),1}$ during the first back-translation loop.

Denoising Bag-of-Word Auto-Encoder (DBAE)

Similarly to both (Lample et al., 2018) and (Wang and Lee, 2018), we also devise a starting model based on a DAE. One major difference is that we use a simple Bag-of-Words (BoW) encoder with fixed pre-trained word embeddings, and a 2-layer GRU decoder. Indeed, we find that a BoW autoencoder trained on the summaries reaches a reconstruction ROUGE-L f-score of nearly 70% on the test set, indicating that word presence information is mostly sufficient to model the summaries. As for the noise model, for each token in the input, we remove it with probability p/2 and add a word drawn uniformly from the summary vocabulary with probability p.

The BoW encoder has two advantages. First, it lacks the other models' bias to keep the word order of the full text in the summary. Secondly, when using the DBAE to predict summaries from the full text, we can weight the input word embeddings by their corpus-level probability of appearing in a summary, forcing the model to pay less attention to words that only appear in \mathcal{D}^F . The Denoising Bag-of-Words Auto-Encoder with input re-weighting is referred to as $f_{F \to S}^{(\mathrm{DBAE}),0}$.

First-Order Word Moments Matching (μ :1)

We also propose an extractive initialization model. Given the same BoW representation as for the DBAE, function $f^{\mu}_{\theta}(s,v)$ predicts the probability that each word v in a full text sequence s is present in the summary. We learn the parameters of f^{μ}_{θ} by marginalizing the output probability of each word over all full text sequences, and matching these first-order moments to the marginal probability of each word's presence in a summary. That is, let \mathcal{V}^S denote the vocabulary of \mathcal{D}^S , then $\forall v \in \mathcal{V}^S$:

$$\mu_v^F = \frac{\sum_{s \in \mathcal{D}^F} \mathbb{1}_{v \in s}}{|\mathcal{D}^F|} \quad \text{and} \quad \mu_v^S = \frac{\sum_{s \in \mathcal{D}^s} \mathbb{1}_{v \in s}}{|\mathcal{D}^S|}$$

We minimize the binary cross-entropy (BCE) between the output and summary moments:

$$\theta^* = \arg\min \sum_{v \in \mathcal{V}^S} \mathrm{BCE}\Big(\frac{\sum_{s \in \mathcal{D}^F} f_{\theta}^{\mu}(s, v)}{|\mathcal{D}^F|}, \mu_v^S\Big)$$

We then define an initial extractive summarization model by applying $f^{\mu}_{\theta^*}(\cdot,\cdot)$ to all words of an input sentence, and keeping the ones whose output probability is greater than some threshold. We refer to this model as $f^{(\mu:1),0}_{F\to S}$.

3.2 Artificial Training Data

We apply the back-translation procedure outlined above in parallel for all three initialization models. For example, $f_{F \to S}^{(\mu:1),0}$ yields the following sequence of models and artificial aligned datasets:

$$f_{F\to S}^{(\mu:1),0} \to \mathcal{D}_{S'\to F}^{(\mu:1),0} \to f_{S\to F}^{(\mu:1),1} \to \mathcal{D}_{F'\to S}^{(\mu:1),1}$$

$$\to f_{F\to S}^{(\mu:1),2} \to f_{S\to F}^{(\mu:1),3} \to \dots$$

Finally, in order to take advantage of the various strengths of each of the initialization models, we also concatenate the artificial training dataset at each odd iteration to train a summarizer, e.g.:

$$f_{F \to S}^{(\mathrm{All}),2} = \mathcal{S} \Big(\mathcal{D}_{F' \to S}^{(\mathrm{Pr\text{-}Thr}),1} \cup \mathcal{D}_{F' \to S}^{(\mathrm{DBAE}),1} \cup \mathcal{D}_{F' \to S}^{(\pmb{\mu}:1),1} \Big)$$

4 Experiments

Data and Model Choices We validate our approach on the Gigaword corpus, which comprises of a training set of 3.8M article headlines (considered to be the full text) and titles (summaries), along with 200K validation pairs, and we report test performance on the same 2K set used in (Rush et al., 2015). Since we want to learn systems from fully unaligned data without giving the model an opportunity to learn an implicit mapping, we also

	R-1	R-2	R-L
Lead-8	21.86	7.66	20.45
PBSMT	24.29	8.65	21.82
Pre-DAE ¹	21.26	5.60	18.89
(Pr-Thr)-0	24.79	8.80	22.46
(DBAE)-0	28.58	6.74	22.72
$(\mu:1)-0$	29.17	8.10	24.71

Table 2: Test ROUGE for trivial baseline and initialization systems. ¹(Wang and Lee, 2018).

further split the training set into 2M examples for which we only use titles, and 1.8M for headlines. All models after the initialization step are implemented as convolutional seq2seq architectures using Fairseq (Ott et al., 2019). Artificial data generation uses top-15 sampling, with a minimum length of 16 for full text and a maximum length of 12 for summaries. ROUGE scores are obtained with an output vocabulary of size 15K and a beam search of size 5 to match (Wang and Lee, 2018).

Initializers Table 2 compares test ROUGE for different initialization models, as well as the trivial Lead-8 baseline which simply copies the first 8 words of the article. We find that simply thresholding on distance during the word alignment step of (Pr-Thr) does slightly better then the full PBSMT system used by Lample et al. (2018). Our BoW denoising auto-encoder with word reweighting also performs significantly better than the full seq2seq DAE initialization used by Wang and Lee (2018) (Pre-DAE). The moments-based initial model (μ :1) scores higher than either of these, with scores already close to the full unsupervised system of Wang and Lee (2018).

In order to investigate the effect of these three different strategies beyond their ROUGE statistics, we show generations of the three corresponding first iteration expanders for a given summary in Table 1. The unsupervised vocabulary alignment in (Pr-Thr) handles vocabulary shift, especially changes in verb tenses (summaries tend to be in the present tense), but maintains the word order and adds very little information. Conversely, the (μ :1) expansion function, which is learned from purely extractive summaries, re-uses most words in the summary without any change and adds some new information. Finally, the autoencoder based (DBAE) significantly increases the sequence length and variety, but also strays from

	Sup.	R-1	R-2	R-L
(Pr-Thr)-2	0	26.17	9.42	23.65
(DBAE)-2	0	28.55	10.24	25.46
$(\mu:1)-2$	0	29.55	9.62	26.10
(All)-2	0	29.80	11.52	27.01
(All)-4	0	30.19	12.36	27.75
(All)-6	0	30.04	12.69	27.64
Advers.	0	28.11	9.97	25.41
REIN-	10K	30.01	11.57	27.61
FORCE ¹	500K	33.33	14.18	30.48
$MASS^2$	100K	29.79	12.75	27.45
FSC ³	500K	30.14	12.05	27.99
Seq2seq ⁴	3.8M	35.30	16.64	32.62

Table 3: Comparison of full systems. The best scores for unsupervised training are bolded. Results from: ¹(Wang and Lee, 2018), ²(Song et al., 2019), ³(Miao and Blunsom, 2016), and ⁴(Nallapati et al., 2016)

the original meaning (more examples in the Appendix). The decoders also seem to learn facts about the world during their training on article text (EDF/GDF is France's public power company).

Full Models Finally, Table 3 compares the summarizers learned at various back-translation iterations to other unsupervised and semi-supervised approaches. Overall, our system outperforms the unsupervised Adversarial-REINFORCE of Wang and Lee (2018) after one back-translation loop, and most semi-supervised systems after the second one, including Song et al. (2019)'s MASS pre-trained sentence encoder and Miao and Blunsom (2016)'s Forced-attention Sentence Compression (FSC), which use 100K and 500K aligned pairs respectively. As far as back-translation approaches are concerned, we note that the model performances are correlated with the initializers' scores reported in Table 2 (iterations 4 and 6 follow the same pattern). In addition, we find that combining data from all three initializers before training a summarizer system at each iteration as described in Section 3.2 performs best, suggesting that the greater variety of artificial full text does help the model learn.

Conclusion In this work, we use the back-translation paradigm for unsupervised training of a summarization system. We find that the model benefits from combining initializers, matching the performance of semi-supervised approaches.

References

- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *TACL*, 5:135–146.
- Kyunghyun Cho, Bart van Merrienboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. In *Proceedings of SSST@EMNLP 2014*, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation, Doha, Qatar, 25 October 2014, pages 103–111.
- Sumit Chopra, Michael Auli, and Alexander M. Rush. 2016. Abstractive sentence summarization with attentive recurrent neural networks. In NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016, pages 93–98.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.
- Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, and Rui Yan. 2018. Style transfer in text: Exploration and evaluation. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018, pages 663–670.
- Sebastian Gehrmann, Yuntian Deng, and Alexander M. Rush. 2018. Bottom-up abstractive summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 November 4, 2018*, pages 4098–4109.
- Edouard Grave, Armand Joulin, and Quentin Berthet. 2018. Unsupervised alignment of embeddings with wasserstein procrustes. *CoRR*, abs/1805.11222.
- Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 448–456.
- Nal Kalchbrenner and Phil Blunsom. 2013. Recurrent continuous translation models. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1700–1709.
- Philipp Koehn, Franz Josef Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Human Language Technology Conference of the North*

- American Chapter of the Association for Computational Linguistics, HLT-NAACL 2003, Edmonton, Canada, May 27 June 1, 2003.
- Guillaume Lample, Myle Ott, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. 2018. Phrase-based & neural unsupervised machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 November 4, 2018*, pages 5039–5049.
- Yishu Miao and Phil Blunsom. 2016. Language as a latent variable: Discrete generative models for sentence compression. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 319–328.
- Ramesh Nallapati, Bowen Zhou, Cícero Nogueira dos Santos, Çaglar Gülçehre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence rnns and beyond. In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany, August 11-12, 2016*, pages 280–290.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. *CoRR*, abs/1904.01038.
- Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 379–389.
- Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 August 4, Volume 1: Long Papers*, pages 1073–1083.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers.*
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. MASS: masked sequence to sequence pre-training for language generation. In Proceedings of the 36th International Conference on Machine Learning, Long Beach, California.
- Sandeep Subramanian, Guillaume Lample, Eric Michael Smith, Ludovic Denoyer, Marc'Aurelio Ranzato, and Y-Lan Boureau. 2018. Multiple-attribute text style transfer. *CoRR*, abs/1811.00552.

- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada, pages 3104—3112.
- Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. 2008. Extracting and composing robust features with denoising autoencoders. In *Machine Learning, Proceedings of the Twenty-Fifth International Conference (ICML 2008), Helsinki, Finland, June 5-9, 2008*, pages 1096–1103.
- Yau-Shian Wang and Hung-yi Lee. 2018. Learning to encode text as human-readable summaries using generative adversarial networks. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 November 4, 2018*, pages 4187–4195.

A Implementation Choices for Initialization and Seq2seq Models

We describe the modeling choices for initialization models (Pr-Thr), (DBAE), and (μ :1). All hyper-parameters for each of these systems are set based on the models' ROUGE-L score on the validation set. Unless otherwise stated, all models use Skipgram FastText¹ word embeddings which are shared across the input and output layers. The dimension 512 embeddings are trained on the concatenation of the full text and summary sequences $\mathcal{D}^F \cup \mathcal{D}^S$. \mathcal{V} is the full vocabulary, and \mathcal{V}^F and \mathcal{V}^S are the vocabularies of \mathcal{D}^F and \mathcal{D}^S respectively. All trained models use the Adam optimizer with learning rate 5e-4. The convolutional seq2seq models use the fconv_iwslt_de_en architecture previded in Fairseq² with pre-trained input and output word embeddings, a vocabulary size of 50K for the full text and of 15K for the summaries. For the expander generations, we collapse contiguous UNK tokens, and cut the sentence at the first full stop even when the model did not generate an EOS token, yielding outputs that are sometimes shorter than 16 words.

Procrustes Thresholded Alignment (Pr-Thr)

For this model, we train two sets of word embeddings on \mathcal{D}^F and \mathcal{D}^S separately, and compute aligned vectors using the FastText implementation of the (Grave et al., 2018) algorithm³. We then map each word in an input sequence to its closest word in \mathcal{V}^S in the aligned space, unless the nearest neighbor is the EOS token or the distance to the nearest neighbor in the aligned space is greater than a threshold η . The output sequence then consists in the first N mapped words in the order of the input sequence. We found that using embeddings of dimension 256, threshold $\eta=0.9$, and maximum output length N=12 yields the best validation ROUGE-L.

We compare (Pr-Thr) to a PBSMT baseline in Table 2. We use the UnsupervisedMT codebase⁴ of (Lample et al., 2018) with the same pre-trained embedding, and also perform a hyper-parameter search over maximum length, which sets N=15.

Denoising Bag-of-Word Auto-Encoder (DBAE)

The DBAE is trained on all sentences in \mathcal{D}^S . The encoder of the DBAE averages the input word embeddings and applies a linear transformation, followed by a Batch Normalization layer (Ioffe and Szegedy, 2015). The decoder is a 2-layer GRU recurrent neural network with hidden dimension 256. The encoder output is concatenated to the initial hidden state of both layers, then projected back down to the hidden dimension.

To use the model for summarization, we perform two changes from the auto-encoding setting. First, we perform a weighted instead of a standard average, where words that are less likely to appear in \mathcal{D}^S than in \mathcal{D}^F are down-weighted (and words that are in \mathcal{V}^F but not in \mathcal{V}^S are dropped). Specifically, given a word $v \in \mathcal{V}^S$, its weight w_v in the summarization weighted BoW encoder is given as:

$$\mu_v^F = \frac{\sum_{s \in \mathcal{D}^F} \mathbb{1}_{v \in s}}{|\mathcal{D}^F|} \text{ and } \mu_v^S = \frac{\sum_{s \in \mathcal{D}^s} \mathbb{1}_{v \in s}}{|\mathcal{D}^S|}$$

$$(5)$$

$$w_v = \max(\frac{\mu_v^S}{\mu_v^F}, 1) \tag{6}$$

Secondly, we implement something like a pointer mechanism by adding λ to the score of each of the input words in the output of the GRU, before the softmax. At test time and when creating artificial data, we decode with beam search and a beam size of size 5, maximum output length N=15, and input word bias $\lambda=2$.

First-Order Word Moments Matching (μ :1)

The moments matching model uses the same encoder as the (DBAE) followed by a linear mapping to the summary vocabulary, followed by a sigmoid layer (the log-score of all words that do not appear in the input is set to -1e6). Unfortunately, computing the output probabilities for all sentences in the corpus before computing the Binary Cross-Entropy is impractical, and so we implement a batched version of the algorithm. Let corpus-level moments μ_v^F and μ_v^S be defined as in Equation 5. Let \mathcal{B}^F be a batch of full text sequences, we define:

$$\hat{\mu}_v^F = \frac{\sum_{s \in \mathcal{B}^F} \mathbb{1}_{v \in s}}{|\mathcal{B}^F|} \text{ and } \hat{\mu}_v^S = \frac{\hat{\mu}_v^F}{\mu_v^F} \cdot \mu_v^S$$
 (7)

¹https://fasttext.cc/

²https://fairseq.readthedocs.io/en/la
test/models.html

³https://github.com/facebookresearch/ fastText/tree/master/alignment

⁴https://github.com/facebookresearch/ UnsupervisedMT/tree/master/PBSMT

For each batch, the algorithm then takes a gradient step for the loss:

$$\hat{\mathcal{L}}(\mathcal{B}^F; \theta) = \sum_{v \in \mathcal{V}^S} \mathrm{BCE}\Big(\frac{\sum_{s \in \mathcal{B}^F} f^{\mu}_{\theta}(s, v)}{|\mathcal{D}^F|}, \hat{\mu}^S_v\Big)$$

The prediction is similar as for the (Pr-Thr) system except that we threshold on $f^{\mu}_{\theta}(s,v)$ rather than on the nearest neighbor distance, with threshold $\eta=0.3$ (the maximum output length is also N=12)

B More Examples of Model Predictions

We present more examples of the expander and summarizer models' outputs in Tables 4, 5, and 6. Table 4 shows more expander generations for all three initial models after one back-translation epoch. They follow the patterns outlined in Section 4, with (DBAE) showing more variety but being less faithful to the input. Table 5 show generations from the expander models at different backtranslation iteration. It is interesting to see that each of the three models slowly overcome their initial limitations: the (DBAE) expander's third version is much more faithful to the input than its first, while the moments-based approach starts using rephrases and modeling vocabulary shift. The Procrustes method seems to benefit less from the successive iterations, but still starts to produce longer outputs. Finally, Table 6 provides summaries produced by the final model. While the model does produce likely summaries, we note that aside from the occasional synonym use or verbal tense change, and even though we do not use an explicit pointer mechanism beyond the standard seq2seq attention, the model's outputs are mostly extractive.

over N,NNN ancient graves found in greek metro dig

(Pr-Thr) over N,NNN ancient graves were found in a greek metro -lrb- UNK -rrb-.

(DBAE) the remains of N,NNN graves on ancient greek island have been found in three ancient graves in the past few days, a senior police officer said on friday.

 $(\mu:1)$ over N,NNN ancient graves have been found in the greek city of alexandria in the northern greek city of salonika in connection with the greek metro and dig deep underground.

ukraine's crimea dreams of union with russia

(Pr-Thr) ukraine 's crimea UNK of the union with russia.

(DBAE) ukraine has signed two agreements with ukraine on forming its european union and ukraine as its membership.

 $(\mu:1)$ ukraine's crimea peninsula dreams of UNK, one of the soviet republic's most UNK country with russia, the itar-tass news agency reported.

malaysian opposition seeks international help to release detainees

(Pr-Thr) the malaysian opposition thursday sought international help to release detainees. the malaysian opposition, news reports said.

(DBAE) malaysian prime minister abdullah ahmad badawi said tuesday that the government's decision to release NNN detainees, a report said wednesday.

 $(\mu:1)$ malaysian opposition parties said tuesday it seeks to "help" the release of detainees.

russia to unify energy transport networks with georgia rebels

(Pr-Thr) russia is to unify energy transport networks with georgia rebels.

(DBAE) russian government leaders met with representatives of the international energy giant said monday that their networks have been trying to unify their areas with energy supplies.

 $(\mu:1)$ russia is to unify its energy and telecommunication networks to cope with georgia's separatist rebels and the government.

eu losing hope of swift solution to treaty crisis

(Pr-Thr) the eu has been losing hope of a UNK solution to the maastricht treaty crisis.

(DBAE) the european union is losing hope it will be a swift solution to the crisis of the eu-lrb- eu-rrb-, hoping that it's in an "urgent" referendum.

 $(\mu:1)$ eu governments have already come under hope of a swift solution to a european union treaty that ended the current financial crisis.

Table 4: More examples of artificial data after the first back-translation iteration.

(Original) malaysia has drafted its first legislation aimed at punishing computer hackers, an official said wednesday.

(Pr-Thr)-1 malaysia has enacted a draft, the first law on a UNK computer hacking.

(Pr-Thr)-3 malaysia has issued a draft of the law on computer hacking.

(Pr-Thr)-5 malaysia has drafted a first law on the computer hacking and internet hacking.

(DBAE)-1 malaysia's parliament friday signed a bill to allow computer users to monitor UNK law.

(DBAE)-3 the country has been submitted to parliament in NNNN passed a bill wednesday in the first reading of the computer system, officials said monday.

(DBAE)-5 malaysia's national defense ministry has drafted a regulation of computer hacking in the country, the prime minister said friday.

 $(\mu : 1)$ -1 malaysia will have drafts the first law on computer hacking.

 $(\mu : 1)$ -3 malaysia has started drafts to be the first law on computer hacking.

(μ : 1)-5 malaysia today presented the nation's first law on computer hacking in the country, news reports said wednesday.

(Title) malaysia drafts first law on computer hacking

Table 5: Evolution of generated full text sequences across iterations.

(Article) chinese permanent representative to the united nations wang guangya on wednesday urged the un and the international community to continue supporting timor-leste.

(Pred) chinese permanent representative urges un to continue supporting timor-leste

(Title) china stresses continued international support for timor-leste

(Article) macedonian president branko crvenkovski will spend orthodox christmas this weekend with the country's troops serving in iraq, his cabinet said thursday.

(Pred) macedonian president to spend orthodox christmas with troops in iraq

(Title) macedonian president to visit troops in iraq

(Article) televangelist pat robertson, it seems, isn't the only one who thinks he can see god's purpose in natural disasters.

(Pred) evangelist pat robertson thinks he can see god's purpose in disasters

(Title) editorial: blaming god for disasters

(Article) the sudanese opposition said here thursday it had killed more than NNN government soldiers in an ambush in the east of the country.

(Pred) sudanese opposition kills N government soldiers in ambush

(Title) sudanese opposition says NNN government troops killed in ambush

Table 6: Example of model predicitons for $f_{F\to S}^{(\text{All},6)}$.