

Understanding Back-Translation at Scale

Sergey Edunov, Myle Ott,
Michael Auli, David Grangier (EMNLP2018)

紹介: 小町研 M1 勝又 智

概要

- NMT で tgt 側の monolingual data を使いたい！
→ back-translation (BT) によって擬似的な対訳データを作成, 学習データに追加する手法がよく用いられる
- この論文は BT の synthetic src について, いくつかの作り方を検討した
- Low-resource の設定で**なければ**, noisy synthetic src による学習データを使用した方が良かった
↑ noisy synthetic src の方が non-noisy より train signal が強いため, という主張
- 他にも, BT による synthetic data と genuine bitext を比較した
- なんか SOTA とった (WMT 2014 En-De, DeepL 比較)

概要

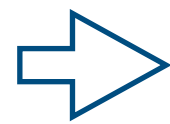
- NMT で tgt 側の monolingual data を使いたい！
→ back-translation (BT) によって擬似的な対訳データを作成, 学習データに追加する手法がよく用いられる
- この論文は **BT の synthetic src** について, いくつかの作り方を検討した

翻訳方向: English → German (自然なデータは緑, 擬似データは橙)

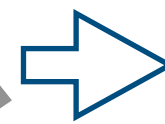
Genuine bitext
(parallel data)

En
Text

De
Text



Train
De-En (BT) Model
with bitext



Make synthetic src
using BT Model
(monolingual data)

De
Text



En
Text

概要

- NMT で tgt 側の monolingual data を学習データとして使う
→ back-translation を用いて synthetic src を生成し、学習データに追加する手法がよく使われる
この **synthetic src** を noisy になるようにした方がいい, というお話
- この論文は **BT の synthetic src** について, いくつかの作り方を検討した

翻訳方向: English → German (自然なデータは緑, 擬似データは橙)

Genuine bitext
(parallel data)

En
Text

De
Text

Train
De-En (BT) Model
with bitext

Make synthetic src
using BT Model
(monolingual data)

De
Text

En
Text

改めて Back-Translation とは

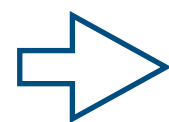
- NMT で tgt 側の monolingual data を使いたい!
→ 逆方向の翻訳モデルを対訳データから学習,
tgt の monolingual data に対して推論を行い, synthetic src を生成する.
この synthetic data を学習データに加えて, 順方向のモデルを学習,
性能向上を目指す. (Sennrich et al., ACL2016)
- 最近は Unsupervised MT でも使われたりしてる. (Lample et al., EMNLP2018)

翻訳方向: English → German (自然なデータは緑, 擬似データは橙)

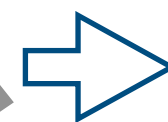
Genuine bitext
(parallel data)

En
Text

De
Text



Train
De-En (BT) Model
with bitext



Make synthetic src
using BT Model
(monolingual data)

De
Text



En
Text

改めて Back-Translation とは

- NMT で **tgt 側の monolingual data** を使いたい!

→ 逆方向の翻訳モデルを対訳データから学習,

tgt 側 Monolingual data は**人手**によるもの. src を生成する.
こ 片側 (monolingual) だけなら適当にクロールすれば を学習,
性 データが得られるので, このデータで何かできたら嬉しい

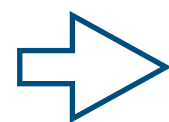
- 最近は Unsupervised MT でも使われたりしてる. (Lample et al., EMNLP2018)

翻訳方向: English → German (自然なデータは**緑**, 擬似データは**橙**)

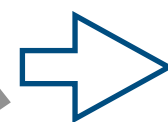
Genuine bitext
(parallel data)

En
Text

De
Text



Train
De-En (BT) Model
with **bitext**



Make synthetic src
using BT Model
(monolingual data)

De
Text



En
Text

改めて Back-Translation とは

- NMT で tgt 側の monolingual data を使いたい!
→ 逆方向の翻訳モデルを対訳データから学習,
tgt の monolingual data に対して推論を行い, synthetic src を生成する.
この synthetic data を学習データに加えて, 順方向のモデルを学習,
性能向上を目指す. (Sennrich et al., ACL2016)
- 最近は Unsupervised MT でも使われたりしてる. (Lample et al., EMNLP2018)

翻訳方向: English → German (自然なデータは緑, 擬似データは橙)

Genuine bitext
(parallel data)

En
Text

De
Text

Train
De-En (BT) Model
with bitext

Make synthetic src
using BT Model
(monolingual data)

De
Text

En
Text

改めて Back-Translation とは

- NMT で tgt 側の monolingual data を使いたい!
→ 逆方向の翻訳モデルを対訳データから学習,
tgt の monolingual data に対して推論を行い, synthetic src を生成する.
この synthetic data を学習データに追加し, 性能向上を目指す. (Sennrich et al., 2015)
- 最近は Unsupervised MT でも使われたりしてる. (Lample et al., 2018)

この BT による synthetic src は
greedy or **beam search** で作るのが通例

翻訳方向: English → German (自然なデータは緑, 擬似データは橙)

Genuine bitext
(parallel data)

En
Text

De
Text

Train
De-En (BT) Model
with bitext

Make synthetic src
using BT Model
(monolingual data)

De
Text

En
Text

改めて

Back-Translation とは

- NMT で tgt 側の monolingual data を使いたい!
→ 逆方向の翻訳モデルを対訳データから学習,
tgt の monolingual data に対して推論を行い, synthetic src を生成する.
この **synthetic data** を学習データに加えて, 順方向のモデルを学習,
性能向上を目指す. (Sennrich et al., ACL2016)

Synthetic data = synthetic src + original tgt

- 最終的に (要するに synthetic data は擬似対訳データ) (Sennrich et al., EMNLP2018)

翻訳方向:

最終的に順方向のモデルは

Gen (parallel)

parallel data (bixtext) + synthetic data で学習するので、
大量のデータで学習することになる。

make synthetic src
using BT Model
(monolingual data)

En
Text

De
Text



De-En (BT) Model
with bixtext



De
Text



En
Text

Synthetic src の作り方 —工夫

- 一般に, BT は synthetic src 作成に greedy or beam search を行う
→ これらの decoding は最も確率の大きいもの (文, 単語) を出力としている
(MAP 推定の出力の近似的なものとみなせる)
- 前提: MAP 推定に近似できる出力は rich translation になりづらい (Ott et al., ICML 2018)
- 仮説: greedy or beam search は出力分布の一番高いところしか見ていないので,
生成される synthetic src は真のデータ分布を適切にカバーしていないのではないかと
→ synthetic src として regular なものを返すのではないかと
- 提案: 出力分布から sampling する or beam 出力に noise を乗せる
- 具体的にはイカの3種の手法を試した
 1. sampling: 出力分布に基づいて sampling (各ステップで)
 2. top-k: 出力分布のうち上位 k token をとり, renormalize, この分布から sampling
 3. beam+noise: beam 出力に対して, 人工的な noise (word drop, blank, shuffle) を乗せる

Synthetic src の作り方

工夫

- 一般に, BT は synthetic src 作成に greedy or beam search を行う
→ これらの decoding は最も確率の大きいもの (文, 単語) を出力としている (MAP 推定 of 出力の近似的なものともみなせる)

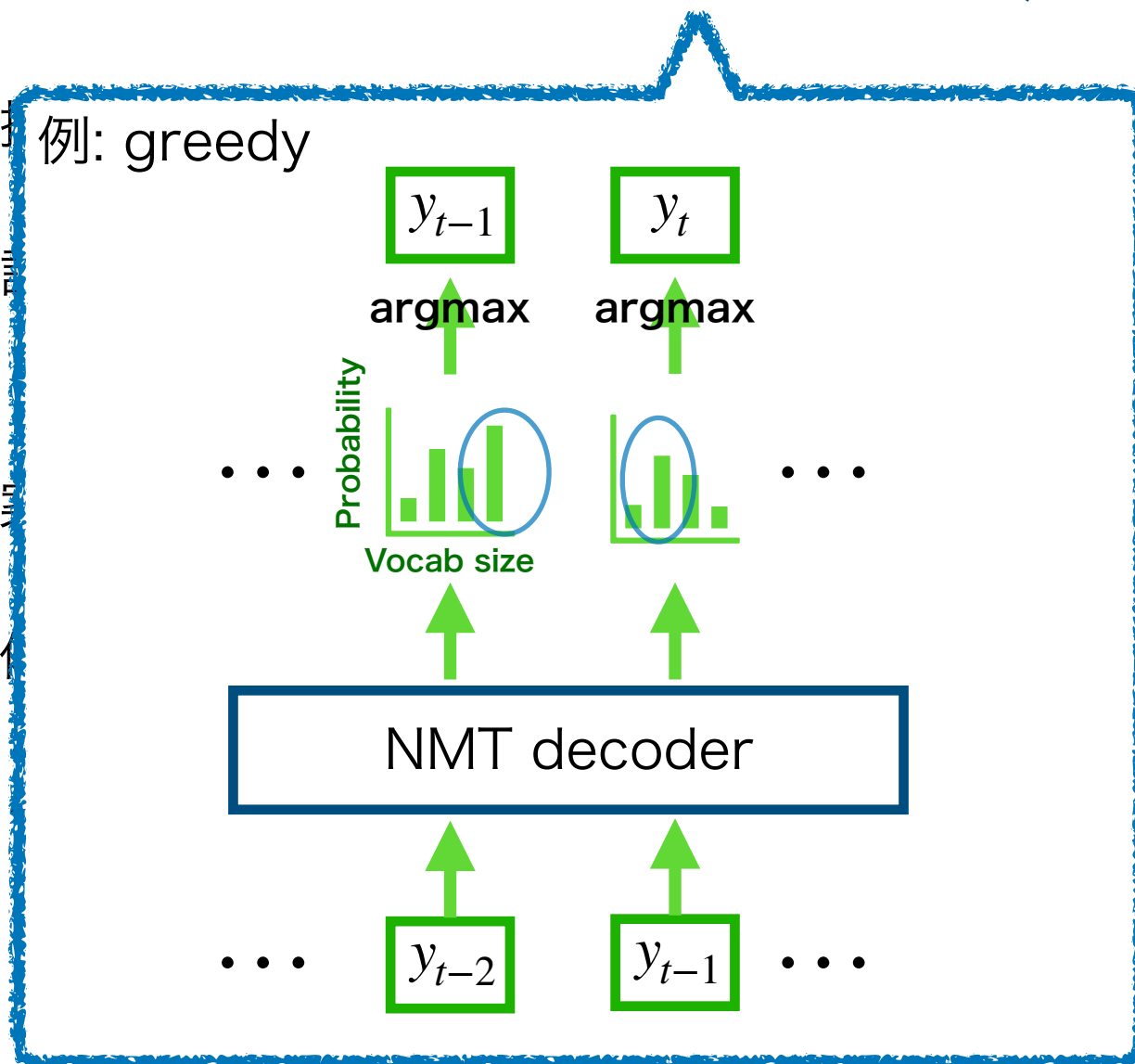
- 前例: greedy

- 仮説

- 提案

- 具体

- 1.
- 2.
- 3.



よりづらい (Ott et al., ICML 2018)

高いところしか見ていないので,
を適切にカバーしていないのではないか?

に noise を乗せる

ステップで)

normalize, この分布から sampling
ise (word drop, blank, shuffle) を乗せる

Synthetic src の作り方 —工夫

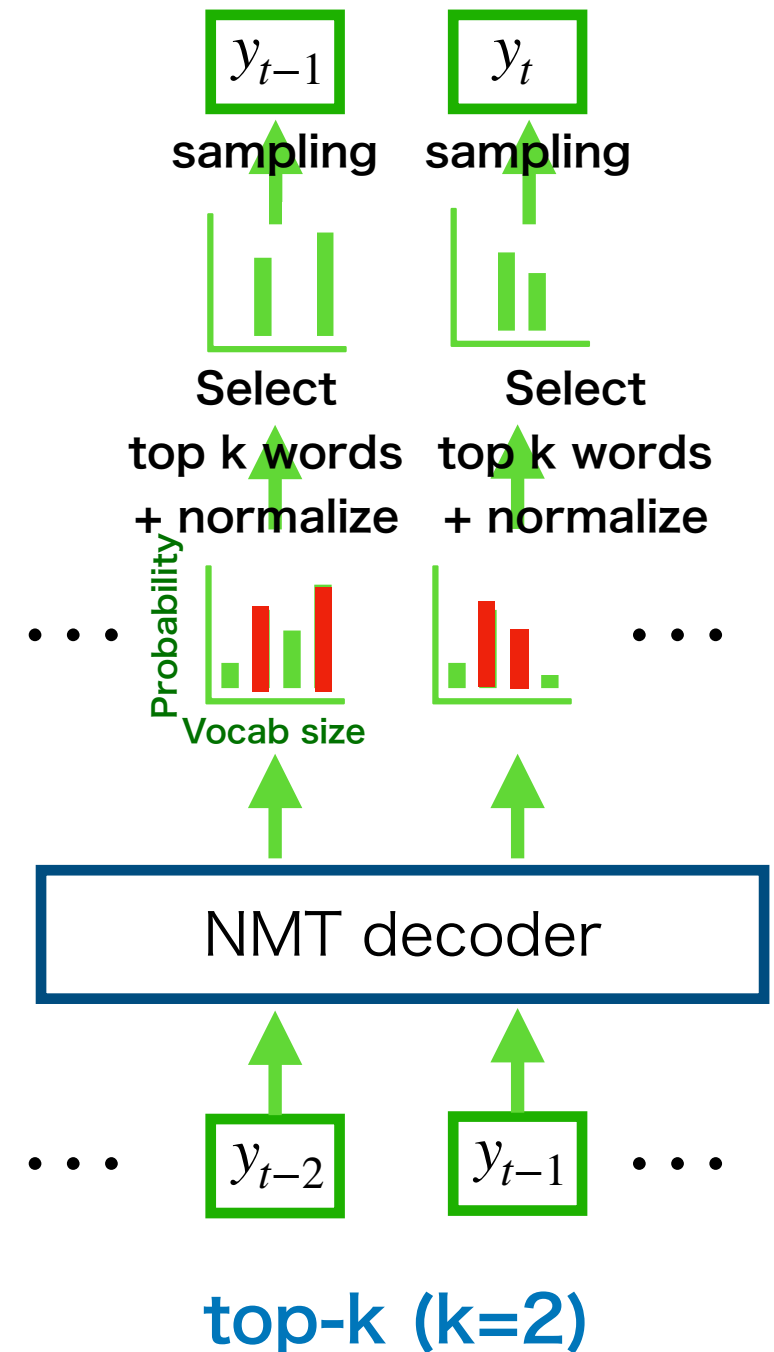
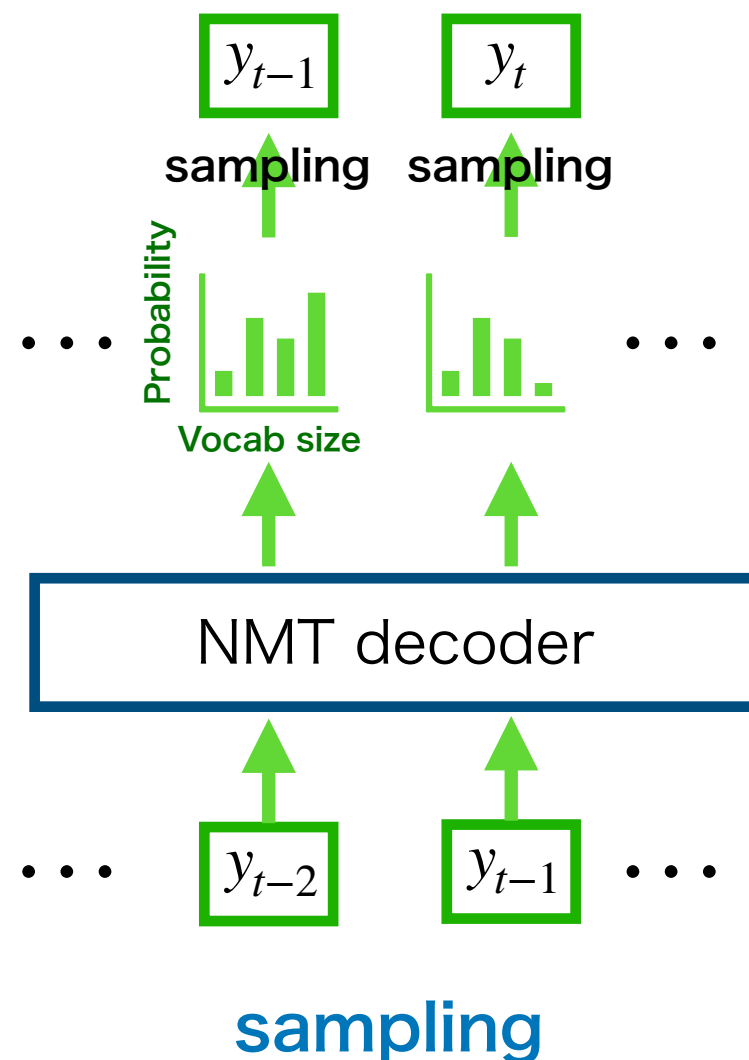
- 一般に, BT は synthetic src 作成に greedy or beam search を行う
→ これらの decoding は最も確率の大きいもの (文, 単語) を出力としている
(MAP 推定の出力の近似的なものとみなせる)
- 前提: MAP 推定による出力は rich translation になりづらい (Ott et al., ICML 2018)
- 仮説: greedy or beam search は**出力分布の一番高いところしか見ていない**ので,
生成される **synthetic src は真のデータ分布を適切にカバーしていない**のではないかと
→ synthetic src として regular なものを返すのではないかと
- 提案: 出力分布から sampling する or beam 出力に noise を乗せる
- 具体的にはイカの3種の手法を試した
1. sampling: 出力分布に基づいて sa
2. top-k: 出力分布のうち上位 k tok
3. beam+noise: beam 出力に対して
beam+noise はそうなるのか…?

気持ちとしては分布の一部を見るのではなく,
もっと多様な感じになるようにしている?
🤔: 離散的にそれをやるのは…どうなんだろう…
る

Syn

- 一般に, BT は s → これらの dec (MAP 推定の出
- 前提: MAP 推定
- 仮説: greedy or 生成される
- 提案: 出力分布が
- 具体的にはイ

↓ 確率が低い単語を選択するかも



1. **sampling**: 出力分布に基づいて sampling (各ステップで)

2. **top-k**: 出力分布のうち上位 k token をとり, renormalize, この分布から sampling

3. beam+noise: beam 出力に対して, 人工的な noise (word drop, blank, shuffle) を乗せる

Synthetic src の作り方

一工夫

- 一般に, BT は synthetic
→ これらの decoding (MAP 推定の出力のよう

- 前提: MAP 推定による出

- 仮説: greedy or beam
生成される synthe

- 提案: 出力分布から sam

- 具体的にはイカの3種の

1. sampling: 出力分布は

2. top-k: 出力分布の

3. **beam+noise**: beam 出力に対して, 人工的な noise (word drop, blank, shuffle) を乗せる

beam 出力が
めちゃくちゃになった文

次の処理を行い, noise を乗せる

- word drop: ある単語を削除する
- word blank: 'BLANK' に置換する
- word shuffle: 順番を入れ替える

beam 探索による一番良い出力

, ICML 2018)

ないので,
ないのではないかな?

から sampling

Synthetic src

この synthetic src を学習に使用する
→ de-noising Auto Encoder

- 一般に, BT は synthetic
→ これらの decoding (MAP 推定の出力の近似)

- 前提: MAP 推定による出力

- 仮説: greedy or beam
生成される synthetic

- 提案: 出力分布から sampling

- 具体的にはイカの3種の

1. sampling: 出力分布から

2. top-k: 出力分布の

3. **beam+noise**: beam 出力に対して, 人工的な noise (word drop, blank, shuffle) を乗せる

beam 出力が
めちゃくちゃになった文

次の処理を行い, noise を乗せる

- word drop: ある単語を削除する
- word blank: 'BLANK' に置換する
- word shuffle: 順番を入れ替える

beam 探索による一番良い出力

, ICML 2018)

ないので,
ないのではないかな?

から sampling

Synthetic src の作り方 —工夫

🤔: sample と beam+noise の出力がなんだかばいことになるのでは？

→ になりました。

(ぶっちゃけ sampling は auto-regressive を考えるともらいことになってもいいと思うが, その辺は語られてない)

source	Diese gegenstzlichen Auffassungen von Fairness liegen nicht nur der politischen Debatte zugrunde.
reference	These competing principles of fairness underlie not only the political debate.
beam	These conflicting interpretations of fairness are not solely based on the political debate.
sample	<i>Mr President</i> , these contradictory interpretations of fairness are not based solely on the political debate.
top10	Those conflicting interpretations of fairness are not solely at the heart of the political debate.
beam+noise	conflicting BLANK interpretations BLANK are of not BLANK based on the political debate.

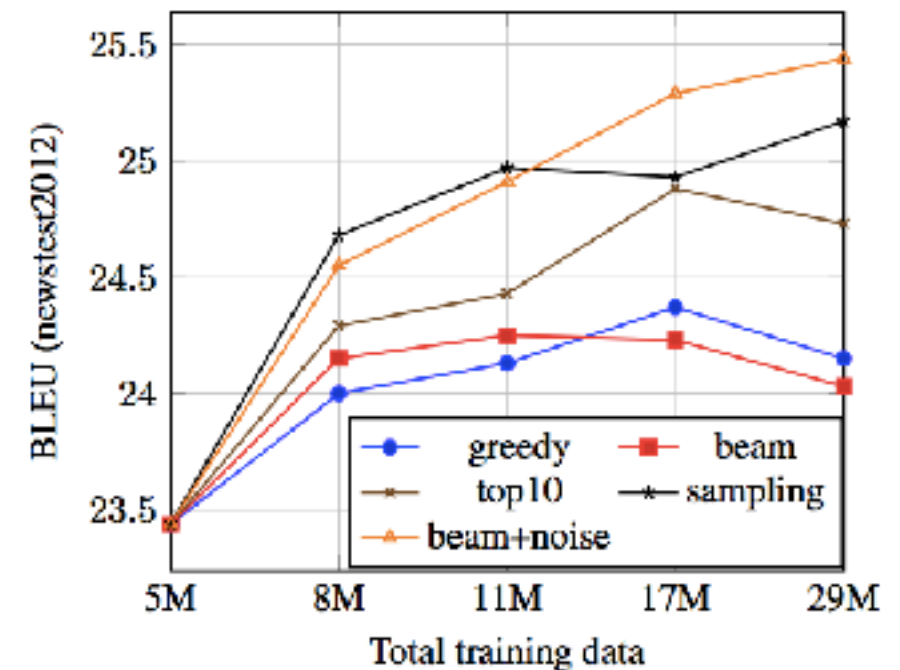
個人的には top-k だと k の幅で多様な出力を出しつつ, そこそこ崩れなくていいのかな, と思っているんですが…

実験設定

- 翻訳方向, data: English-German, WMT'18 ← main
data size: [parallel] 5.18M sents; [monolingual] 226M sents (max)
dev: newstest2012
test: newstest2013-2017
学習データ (parallel) の 52K sents を held-out set としてる ← 擬似データを入れた時の更新回数を決めてる
- BPE (joint src and tgt): 35K types ← sampling ですごい出力が出てきそう
- 翻訳方向, data: English-French, WMT'14 ← SOTA が言いたい
data size: [parallel] 35.7M sents; [monolingual] 31M
dev: newstest2012
test: newstest2013-2015
- 評価: tokenized BLEU と detokenized BLEU どちらもやった
- NMT モデル: Big Transformer (fairseq)
beam size: 5
- Top-k: k= {5, 10, 20, 50} でやって, どれも似たような結果になったので **k=10** を報告

Synthetic data の作り方を 色々やった結果

5M が bitext only, そこに擬似データを足してる. →



↓ bitext 5.2M sents + synthetic sents 24M

	news2013	news2014	news2015	news2016	news2017	Average
bitext	27.84	30.88	31.82	34.98	29.46	31.00
+ beam	<u>27.82</u>	32.33	32.20	35.43	31.11	31.78
+ greedy	<u>27.67</u>	32.55	32.57	35.74	31.25	31.96
+ top10	28.25	33.94	34.00	36.45	32.08	32.94
+ sampling	28.81	<u>34.46</u>	<u>34.87</u>	37.08	<u>32.35</u>	<u>33.51</u>
+ beam+noise	<u>29.28</u>	33.53	33.79	<u>37.89</u>	32.66	33.43

Table 1: Tokenized BLEU on various test sets of WMT English-German when adding 24M synthetic sentence pairs obtained by various generation methods to a 5.2M sentence-pair bitext (cf. Figure 1).

Synthetic data の作り方を 色々やった結果

sampling と beam+noise がいい感じだった。

top10 はちょっといい。

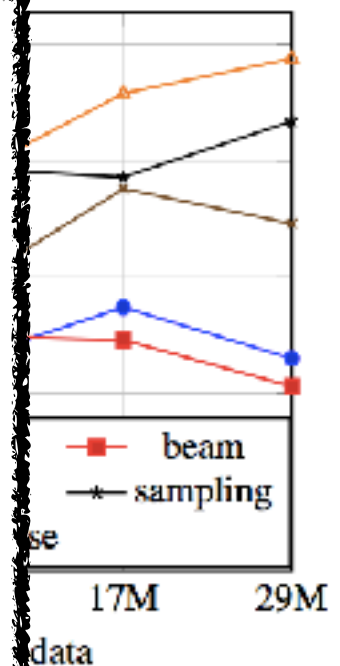
→ なんでだろう？

noisy src sentence は学習が hard. (denoising AE 並感)

Sampling 出力は argmax 出力より分布をよく表現できている。

→ どちらも **training signal** が beam や greedy より richer なのは？

→ その辺りを分析する。



	news2013	news2014	news2015	news2016	news2017	Average
bitext	27.84	30.88	31.82	34.98	29.46	31.00
+ beam	<u>27.82</u>	32.33	32.20	35.43	31.11	31.78
+ greedy	<u>27.67</u>	32.55	32.57	35.74	31.25	31.96
+ top10	28.25	33.94	34.00	36.45	32.08	32.94
+ sampling	28.81	<u>34.46</u>	<u>34.87</u>	37.08	<u>32.35</u>	<u>33.51</u>
+ beam+noise	<u>29.28</u>	33.53	33.79	<u>37.89</u>	32.66	33.43

Table 1: Tokenized BLEU on various test sets of WMT English-German when adding 24M synthetic sentence pairs obtained by various generation methods to a 5.2M sentence-pair bitext (cf. Figure 1).

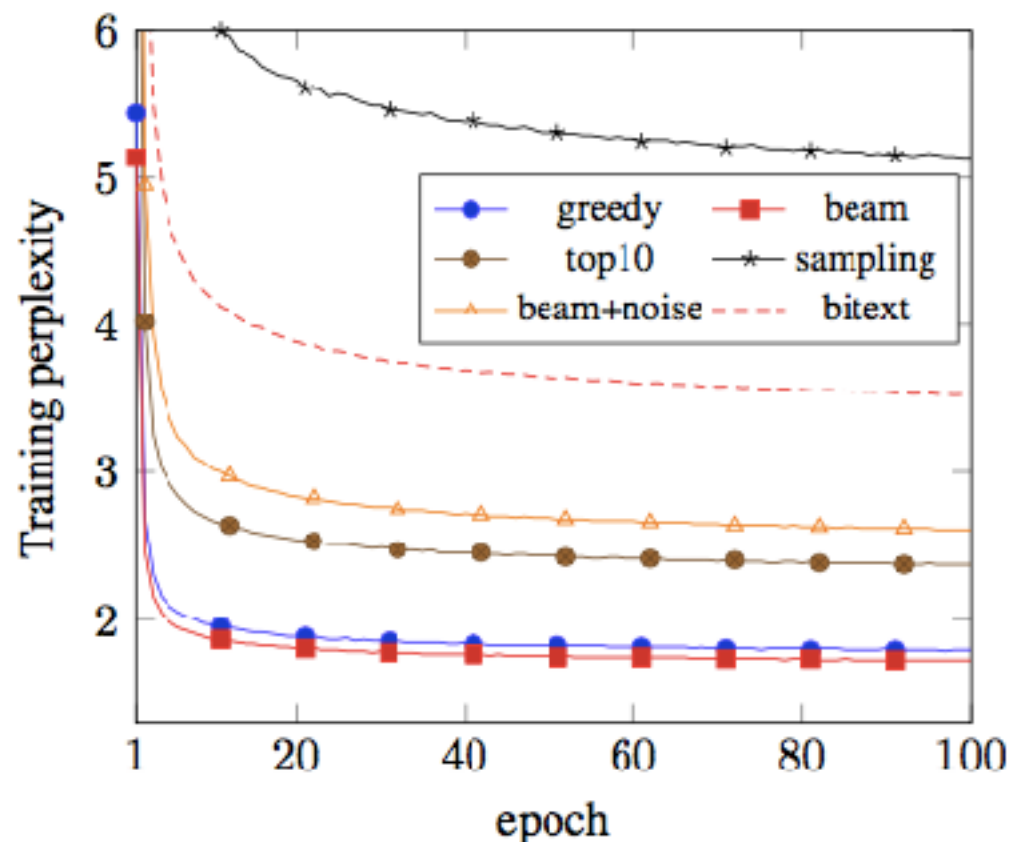
richer training signal なのかを 調べて見た with PPL

training loss (ppl) をみてみた (synthetic data: 24M sents)

各 epoch 終了時に training data の 500K 文対に関する ppl

- bitext only は bitext から 500K;
- BT 有りは synthetic data から 500K (tgt 側として同一のものを使用)

→ ppl が即下がる場合は学習がちょろい, ということになる



1. Sampling や beam+noise は ppl が大きい

→ richer training signal, ということらしい

🤔: そもそも何をもって **rich** training signal と言うのか. ppl bigger → richer は言える…?

2. sampling や noise によって多様な src を見る

→ model は reordering や substitution に robust になる (たとえ擬似 noise が現実には存在しなくても)

🤔: 部分的には有りそうだけでも…それにしても上がりすぎでは?

richer training signal なのかを 調べて見た with 言語モデル

- めっちゃ regular な 学習 data は Language Model で当てやすいのでは
(つまり多様な感じが regular な感じかを評価する)
- 要するに
LM を bitext (tgt) で学習,
分野の同じ bitext で BT を学習,
この BT によって生成された synthetic src を LM で測る
low ppl → regular synthetic src
- 実際には, bitext から
LM 学習の学習に 4.1M sents (src 側) (5-gram; Kneser-Ney)
BT モデル学習に 640K sents,
450K に対して BT, 得られた synthetic src に対して LM の ppl を計測

richer training signal なのかを 調べて見た with 言語モデル

- めっちゃ regular な 学習 data は Language Model で当てやすいのでは
(つまり多様な感じか regular な感じかを評価する)

- 要 - この結果から, sampling や beam+noise の方が
LM rich な出力を出している, と言える

分 - 単純な beam では 変動性が不足してて
こ noise とかよりも training signal が弱い
low みたいなことが言えそうですね.

- bit ということらしいが...

LM

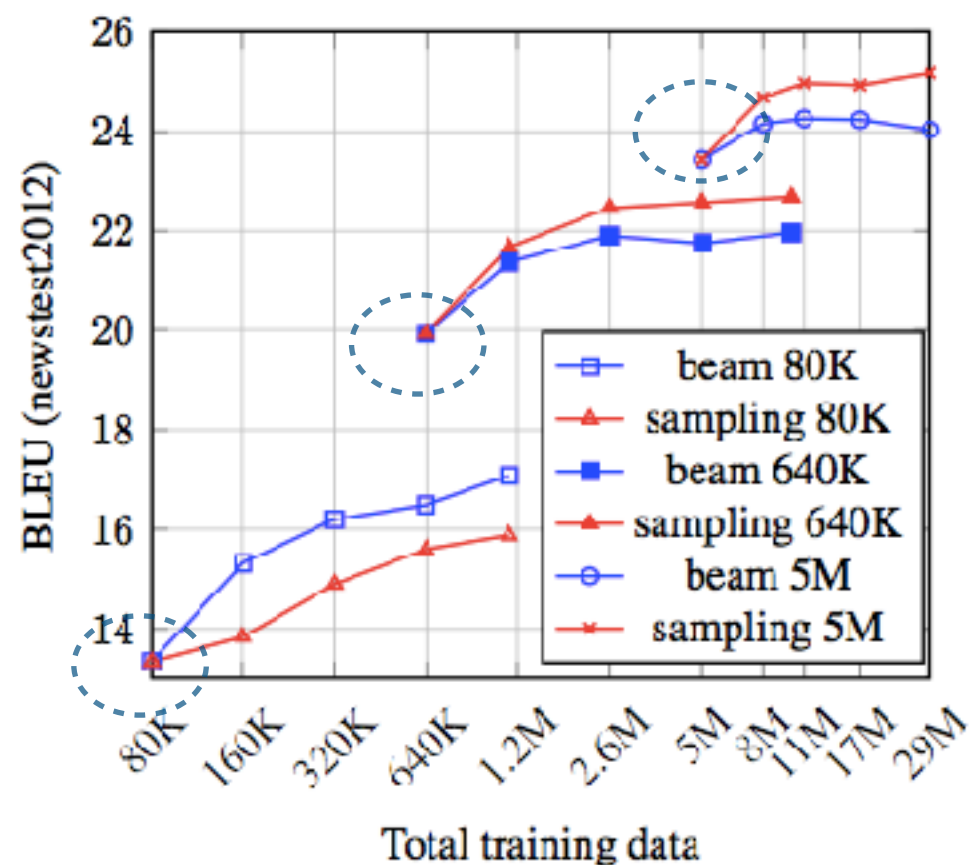
BT

Perplexity	
human data	75.34
beam	72.42
sampling	500.17
top10	87.15
beam+noise	2823.73

450K に対して BT, 得られた synthetic src に対して LM の ppl を計測

Low-resource で効くとめっちゃ 嬉しくないですか？

- 今までは bitext として 5M ぐらい使ってきた
→ 減らした時に, この noise を含んだ synthetic src は有効か？
- bitext を 80K 使う設定, 640K 使う設定でやって見た
BT (De-En) は 80K で 13.5 BLEU, 640K で 24.3 BLEU, 5M で 28.3 BLEU



← 点線の丸で囲んだ点が bitext only
そこから synthetic src を足してる形

結果: low-resource だと sampling で増やす
よりは beam で増やした方が良さそう
ていうのも, きちんとしたコーパスが小さいと
そもそも noise に対して脆い感じになるから.

Synthetic data の domain 事情

- 本物 vs 偽物, Domain とかの影響調査する

BT は sampling を行う

- Bitext から 640K sents pair を使って BT を学習する
さらにこの 640K sents pair と
 1. 残りの parallel data を使用して順方向を学習 (normal bitext)
 2. 残りの parallel data の tgt に対して BT による推論を行い, その synthetic data を使用して順方向を学習 (BT-bitext)
 3. Newscrawl data に対して BT による推論を行い, その synthetic data を使用して順方向を学習 (BT-news)
- 1 と 2 の比較で 本物 vs 偽物 を,
2 と 3 の比較で, domain に関する影響調査を行う.
- 評価 set: newstest2012 (newswire) と WMT の held-out set (mix)

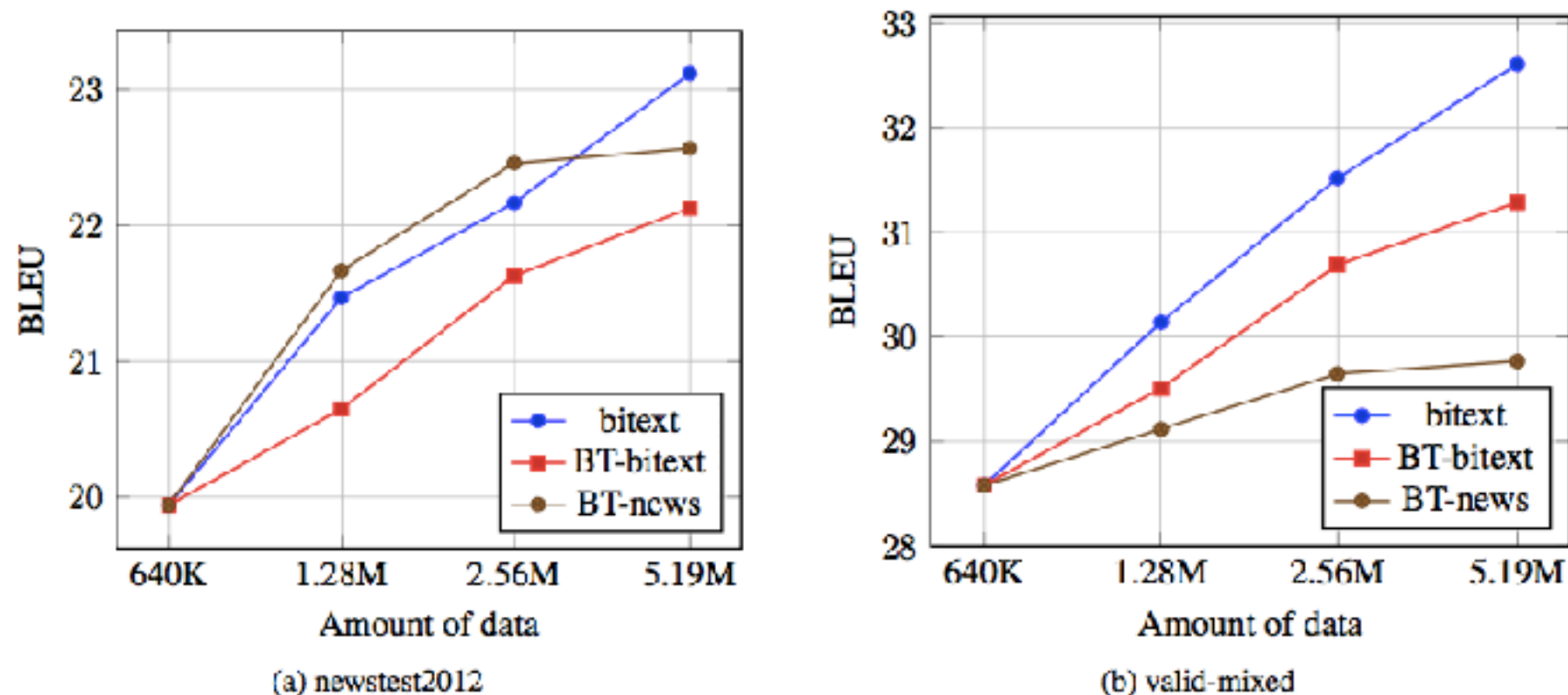


Figure 4: Accuracy on (a) newstest2012 and (b) a mixed domain valid set when growing a 640K bitext corpus with (i) real parallel data (bitext), (ii) a back-translated version of the target side of the bitext (BT-bitext), (iii) or back-translated newscrawl data (BT-news).

1. 本物 vs 偽物

a → 本物 (bitext) が 2.6 BLEU up; 偽物 (BT-bitext) が 2.2 BLEU up
→ 本物に対して83%まで達成!

2. Domain 調査

a, b の BT-news, BT-bitext の改善度合いをみる → 揃っている方が良い
a から十分量の同一 domain monolingual があれば, もともと mix で学習 (bitext) した場合よりもいい感じ (domain adaptation 的なことはまあできる)

bitext の方を upsampling して見た

- Synthetic data と data size を揃えるような気持ち

- 具体的にはもともと bitext 5M だったものを n 倍する感じ

例: $n = 4$

→ bitext 20M, synthetic 24M

- Fig 5 が n を 1, 2, 4, 8 にした時の結果
Synthetic data は 24M 固定

- Greedy や beam は n をあげると良いが, noise 系はそうでもない

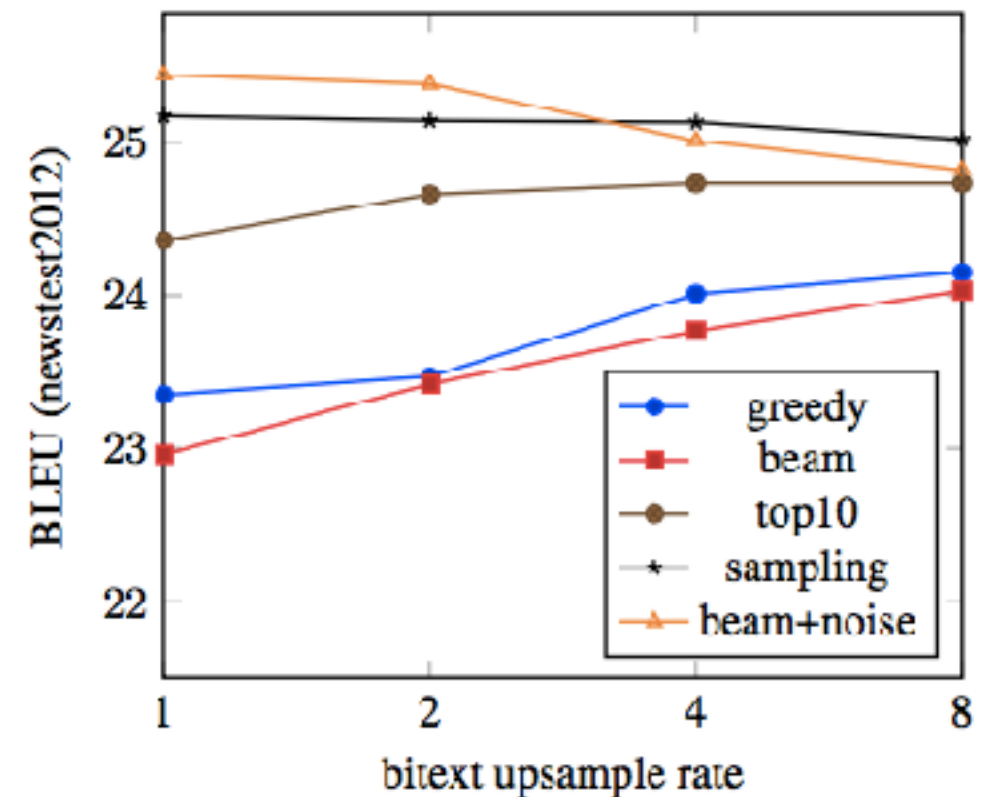


Figure 5: Accuracy when changing the rate at which the bitext is upsampled during training. Rates larger than one mean that the bitext is observed more often than actually present in the combined bitext and synthetic training corpus.

SOTA とった話

- monolingual data として 226M sents 突っ込んだ
→ 取りました

	En-De	En-Fr
a. Gehring et al. (2017)	25.2	40.5
b. Vaswani et al. (2017)	28.4	41.0
c. Ahmed et al. (2017)	28.9	41.4
d. Shaw et al. (2018)	29.2	41.5
DeepL	33.3	45.9
Our result	35.0	45.6
<i>detok. sacreBLEU</i> ³	33.8	43.8

Table 6: BLEU on newstest2014 for WMT English-German (En-De) and English-French (En-Fr). The first four results use only WMT bitext (WMT'14, except for b, c, d in En-De which train on WMT'16). DeepL uses proprietary high-quality bitext and our result relies on back-translation with 226M newscrawl sentences for En-De and 31M for En-Fr. We also show detokenized BLEU (SacreBLEU).

まとめ

- Back translation による出力を noisy な感じにすることで性能が上がった.
→ 単純な従来 of 出力と比べて, richer training signal を出してるからでしょう. (ただし low-resource は除く)
- SOTA とった
- 🤔: 色々雑だけでもとりあえず大量データがあるなら離散的 noise (src 側のみ) は割と有効?
最初は n-best 出力, normalize, sample sentence みたいにすると思ってた勢

参考文献

- Improving Neural Machine Translation Models with Monolingual Data. Sennrich et al., ACL2016
- Phrase-Based & Neural Unsupervised Machine Translation. Lample et al., EMNLP2018
- Analyzing Uncertainty in Neural Machine Translation. Ott et al., ICML 2018
- GEC で似た話
Noising and Denoising Natural Language: Diverse Backtranslation for Grammar Correction. Xie et al., NAACL 2018