Understanding Knowledge Distillation in Neural Sequence Generation

Jiatao Gu January 22, 2020

facebook Artificial Intelligence Research



Jiatao Gu

Ph.D.

- New York, US
- Facebook AI Research

About Me

I am currently a research scientist at the <u>Facebook AI Research</u> in New York City. My general research interests lie in applying deep learning approaches to natual language processing (NLP) problems. In particular, I am interested in building an efficient, effective and reliable neural machine translation (NMT) system for human languages.

I obtained my Ph.D. degree at the department of Electrical and Electronic Engineering, University of Hong Kong in 2018 and I was supervised by <u>Prof. Victor O.K. Li</u>. I spent a wonderful time visiting the <u>CILVR Lab</u>, New York University working with <u>Prof. Kyunghyun Cho</u>. Before that, I obtained my Bachelor's degree at the Electronic Engineering Department, Tsinghua University in 2014 with the guidance of <u>Prof. Ji Wu</u>.

Towards Better Understanding & Interoperability of NMT

Low-Resource and Multilingual NMT

Flexible Representation and Efficient Decoding for NMT

End-to-end Simultaneous Speech Translation

facebook Artificial Intelligence Research





Low-Resource and Multilingual Neural Machine Translation

- Improved Zero-shot NMT (Gu et al. 2019, ACL 2019)
- Multilingual NMT with Byte-level subwords (Wang et al. 2019, AAAI 2020) •
- **TACL 2020)**



The Source-Target Domain Mismatch Problem in NMT (Shen et al. 2019, submitted to

Multilingual Denoising Pre-training for NMT (arxiv today)



して証明と証拠を求めましょう	Ask_questions,demandproof,demandevidence.
3 AA E5 95 8F E3 81 97 E3 81 A6 E2 96 81 E8 A8 BC E6 98 8E A8 E8 A8 BC E6 8B A0 E3 82 92 E6 B1 82 E3 82 81 E3 81 BE 97 E3 82 87 E3 81 86	41 73 6B E2 96 81 71 75 65 73 74 69 6F 6E 73 2C E2 96 81 64 65 6D 61 6E 64 E2 96 81 70 72 6F 6F 66 2C E2 96 81 64 65 6D 61 6E 64 E2 96 81 65 76 69 64 65 6E 63 65 2E
3 AA E595 8F しE381 A6 <mark>E8 A8 BC 明 E381 A8 E8 A8 BC E6</mark> D をE6 B1 82 めE381 BE しょう	As kquest ions ,dem andpro of ,dem andev idence .
3 AA 問 しE381 A6 <mark>E8 A8BC 明 E381 A8 E8 A8BC E68B A0</mark> を 1 82 めE381 BE しょう	A s kqu est ion s ,d em andpro o f ,d em ande v id ence .
3 AA 問 しE381 A6 <u>E8 A8BC 明E381 A8 E8 A8BC</u> 拠 をE6 B1 E381 BE しょう	A s kquest ions ,d em andpro of ,d em andev id ence .
3 AA問 しE381 A6 <u></u> E8 A8BC 明E381 A8 E8 A8BC 拠 をE6 B1 E381 BE しょう	As kquestions ,demandpro of ,demandevidence .
3 AA問 しE381 A6 <u></u> E8 A8BC 明E381 A8 E8 A8BC 拠 をE6 B1 E381 BE しょう	As kquestions ,demandproof ,demandevidence .
3 AA問しE381 A6 <u>E8 A8BC 明E381 A8 E8 A8BC</u> 拠 をE6 B1 82 81 BE しょう	As kquestions ,demandproof ,demandevidence .
して_証明と証拠を求めましょう	Ask_questions,_demand_proof,_demand_ evidence.
して 証明 と 証拠 を求 め ましょう	As kquestions ,demandpro of ,demandevidence .
して証明 と 証拠 を求め ましょう	As kquestions ,demandproof ,demandevidence .



Flexible Representation and Efficient Decoding for NMT

- Insertion-based Generation (Gu et al. 2019, TACL 2019)
- Levenshtein Transformer (Gu et al. 2019, NeurIPS 2019)

coder depth

De

- Generation with Adaptive Computational Time (Elbayad et al. 2019, ICLR 2020)
- Parallel Machine Translation with Disentangled Context Transformer (Kasai et al. 2020, to ICML2020)







facebook Artificial Intelligence Research





End-to-end Simultaneous Speech Translation

...



End-to-End AST with Indirect Training Data (Pino et al. 2019, IWSLT 2019) Simultaneous Speech Translation (Ma et al. 2019, submitted ICLR2020) Multilingual Speech Translation (submitted to LREC2020)

	IWSLT2020		\heartsuit	2	Search!	Q
	Simultaneous_translati `	Simultaneous Speech Translation				
	 Conference Evaluation Important Dates Organizers Past editions Trace: start evaluation simultaneous_translation 	Task Description Simultaneous machine translation has become an increasingly popular top speech translation (SST) enables interesting applications such as subtitle translation. The goal of this task is to examine systems for translating audic language with consideration of both translation quality and latency, and the research community on this direction. We encourage participants to submit systems either based on cascaded (A participants will be evaluated on translating TED talks from English into Ge enter: Text-to-Text (T2T-SST): participants will be asked to translate the goal Text (S2T-SST): participants need to directly translate the audio speech interer both tracks when possible. Evaluating a simultaneous system is not trivial as we cannot release the tess participants will be required to implement specific APIs to read the input a systems as a Docker image where we will evaluate on our own environment which will also serve as the baseline system. The system's performance will be evaluated in two folds: The translation qual BLEU, TER, and ChrF. The translation latency: we will make use of the recent translation including average proportion (AP), average lagging (AL) and diff we will report timestamps for informational purposes. We will provide the with the Docker example.	ic in recent translation o speech in the ultimate ASR + MT) rman. They ound-truth to text in react t data as o nd write the t. We will p uality: we we tly develop ferentiable example of	t years. In a for a live a one lang goal is to o or end-to y will be gi transcript al-time. W offline trans al-time translati provide an will use mu ped metric a verage l f computir	particular, simulta event or real-time uage into text in the foster advances fr o-end approaches ven two parallel tr s in real-time. Spe- e encourage partice slation tasks do. In on, and upload the example implement ltiple standard me s for simultaneous agging (DAL). In a ng these metrics to	aneous video call he target form the . This year, racks to ech-to- cipants to astead, eir entation etrics: a machine ddition, ogether
		Contacts Chair: Jiatao Gu (Facebook, USA) Discussion: ⊠iwslt-evaluation-campaign@googlegroups.com				
Welco partici	me to pate!	• Jiatao Gu (Facebook)				

- Jiatao Gu (Facebook)
- Juan Pino (Facebook)



Low-Resource and Multilingual NMT

Flexible Representation and Efficient Decoding for NMT

End-to-end Simultaneous Speech Translation

facebook Artificial Intelligence Research



Towards Better Understanding & Interoperability of NMT

7

Sequence-Level Knowledge Distillation

facebook Artificial Intelligence Research

Neural Machine Translation

Model:

- Encoder/Decoder \rightarrow RNNs (before 2018) / Transformers (now)
- We model the conditional probability in *autoregressive* factorization. •



 $\mathbf{L} \mathbf{L}_{t=1}$

 $P(Y|X) = | P(y_t|y_{1...t-1}, x_{1...T'}; \theta)$

it ties	
x	
<u>™</u>)	
d	
ad m	N×
d ad on	
<u>_</u>)	
	sitional coding
t ing	0

9

Neural Machine Translation

Data:

- Parallel corpus (sentence-level aligned)
- Weakly supervised corpus
- Monolingual corpus \rightarrow semi-supervised/unsupervised learning •





 $P(Y|X) = \left[P(y_t|y_{1...t-1}, x_{1...T'}; \theta) \right]$

it ties	
x	
<u>™</u>)	
d	
mr.	
ad m	N×
d ad on	
	sitional coding
t ing	-



Neural Machine Translation

• Training: maximum likelihood training

$$J^{ML}(\theta) = \sum_{t=1}^{T} \log P(y_t | y_{1...t-1}, x_{1...T'}; \theta)$$

Decoding: greedy / beam-search

$$\hat{y}_{t} = \underset{y}{\operatorname{argmax}} \log P(y_{t} | \hat{y}_{t-1,\dots,1}, x_{1\dots T'}; \theta) \qquad \text{Attention}$$

$$X = x_{1}, x_{2}, \dots, x_{T'} \longrightarrow \text{Encoder}$$

 $P(Y|X) = P(y_t|y_{1...t-1}, x_{1...T'}; \theta)$ $\mathbf{L} \mathbf{L}_{t=1}$



it ties	
x	
<u>™</u>)	
d	
mr.	
ad m	N×
d ad on	
	sitional coding
t ing	-



Knowledge Distillation

- Knowledge distillation (Liang et al., 2008; Hinton et al., 2015) was originally proposed for training a weaker student classifier on the targets predicted from a stronger teacher model.
- A typical approach is using the label probabilities produced by the teacher as "soft targets" (dark knowledge)

 $q_i = \frac{\exp(z_i/\tau)}{\sum_i \exp(z_i/\tau)}$

• In the context of sequence generation, Kim & Rush (2016) extend this idea using "hard targets" from a teacher generation model. More precisely, $q(t|x) \approx \mathbb{I}\{t = \operatorname{argmax}_{t \in T} q(t|x)\}$:

$$egin{aligned} \mathcal{L}_{ ext{seq-KD}} &= -\mathbb{E}_{m{x} \sim ext{data}} \sum_{m{t} \in \mathcal{T}} q(m{t} | m{x}) \log p(m{t} | m{x}) \ &pprox & -\mathbb{E}_{m{x} \sim ext{data}}, \hat{m{y}} = rg\max_{m{t} \in \mathcal{T}} q(m{t} | m{x}) \left[\log p(m{t} = \hat{m{y}} | m{x})
ight] \end{aligned}$$

12

Sequence-level Knowledge Distillation



A Teacher-Student Framework in Three Steps:

(1) Train a teacher model with golden targets.

(2) Generate new targets with the pretrained teacher.

(3) Train the student model with the generated targets.

13

Sequence-level Knowledge Distillation



Questions:

- generation?

(1) How to choose the teacher/student models?

(2) What kind of data can we use for distillation?

(3) In fact, why and how does distillation work in

14

Understanding Knowledge Distillation in Non-autoregressive Machine Translation

w/ Chunting Zhou and Graham Neubig

facebook Artificial Intelligence Research

ICLR2020





Non-autoregressive Neural Machine Translation



(Figure from Gu et.al, 2017)

- Strong: Autoregressive model (e.g. Transformers) can in theory model any arbitrary distribution of sequences.
- **Slow:** we need to predict one word and a time during inference.

in parallel:

- Fast: An alternative solution where we predict all the target tokens in parallel which is favorable for parallelism.
- independent.

facebook Artificial Intelligence Research Standard NMT systems are *autoregressive (AT model)*:

$$P(Y|X) = \prod_{t=1}^{T} P(y_t|y_{1:t-1}, x_{1:T'})$$

Non-autoregressive Translation (NAT model) predicts sequence generation

Weak: It is harmful to assume all the output tokens are completely

$$P(Y|X) = \prod_{t=1}^{T} P(y_t | x_{1:T'})$$

16

Levenshtein Transformer

Our major contribution is still at the "modeling side":

- Iterative-based Parallel Refinement (a Markov Decision Process)
- For each iteration, we extend the model by considering "deletion" and "insertion" as the basic operations.
- Both insertion and deletion operations are "non-autoregressive"!



$$\pi(\boldsymbol{a}|\boldsymbol{y}) = \prod_{d_i \in \boldsymbol{d}} \pi^{\text{del}}(d_i|i, \boldsymbol{y}) \cdot \prod_{p_i \in \boldsymbol{p}} \pi^{\text{plh}}(p_i|i, \boldsymbol{y}') \cdot \prod_{t_i \in \boldsymbol{t}} \pi^{\text{tok}}(t_i|i, \boldsymbol{y}'')$$

Model



- **Encoder-decoder attention** was omitted.
- The parameters of three passes can be shared.
- We also propose to "Early **Exit**" which attaches the classifiers to an intermediate block instead of the last to save computation.



Non-autoregressive Neural Machine Translation

In practice, it is always helpful to obtain some forms of intermedia representation Z to capture the ignored dependency between output tokens in NAT.

For instance,



...

...



<s> a cat sat on the mat </s

(Figure from Gu et.al, 2019)

facebook Artificial Intelligence Research

$$P(Y|X) = \sum_{Z} P(Z|x_{1:T'}) \cdot \prod_{t=1}^{T} P(y_t|Z, x_{1:T'})$$

Two types of NAT-based models are often considered: • Z as standard discrete/continuous latent variables (VAE-based NAT) -- <u>https://arxiv.org/abs/1803.03382</u> -- <u>https://arxiv.org/abs/1909.02480</u>

• Z as intermedia partial generation (Refinement-based NAT) -- <u>https://www.aclweb.org/anthology/D18-1149/</u> -- <u>https://papers.nips.cc/paper/9297-levenshtein-transformer.pdf</u>



Knowledge Distillation for NAT



As one of the most successful tricks, KD has been used in *almost* all existing NAT models.

• Typically, the student is our targeted NAT model, while we choose the teacher an autoregressive model (AT).

• As discussed earlier, we can assume "teacher" is much stronger than the student to model the data.

Both teacher and student models are trained on the same source sentences.



Knowledge Distillation for NAT



Here is the example performance w/ and w/o distillation for NAT models.

- Test set BLEU on WMT14 English-German (En-De)
 - All three models distilled from the same AT Transformer with BLEU score of 27.13 on WMT En-De.

	w/o distillation	w/ distillation
NAT (Gu et al, 2017)	11.4	19.5 (+8.1)
q (Ma et al, 2019)	18.6	21.7 (+3.1)
Gu et al, 2019)	25.2	26.9 (+1.7)

How does knowledge distillation improve NAT models so much?



Multi-modality Problem

The original NAT paper (Gu et al, 2017) argues the fundamental issue for nonautoregressive models as the multi-modality problem in the data:

For example:



Thank you

Our assumption is that distillation helps to reduce the multimodality in the data.





Case study on Toy Data



When things are unclear and too difficult to explain in sequence generation (e.g. machine translation tasks), it is always a good idea to look at some toy cases.

- language to be output.

We manually created the multi-modality (language id) in the data.

We create a synthetic dataset compared with three language pairs -- English-German (En-De), English-French (En-Fr) and English-Spanish (En-Es) – from the Europarl corpus. We make sure every English sentence will be aligned to ALL three languages, and no language ID was specified.

• We train both AT and NAT models directly on this synthetic dataset. During inference time, we input the English sentence without telling the model which

23



facebook Artificial Intelligence Research

$$p(l_i | \mathbf{y}) \approx \frac{1}{T} \sum_{t=1}^T p(l_i | y_t) = \frac{1}{T} \sum_{t=1}^T \frac{p(y_t | l_i) p(l_i)}{\sum_k p(y_t | l_k) p(l_k)}$$

but distillation is a more systematic way of mode selection.

Case study on Toy Data

Inspired from the visualization on toy data, we propose to use "data uncertainty" to measure the multi-modality (complexity) for **general purpose**.

For simplicity, the data uncertainty is calculated by fitting an alignment model (we use fast-align) and compute the average of **token-level conditional entropy**.





The corpus level complexity is a simple average of the token-level conditional entropy over the vocabulary.

C(d) =



facebook Artificial Intelligence Research

$$\sum_{\substack{\in \mathcal{Y} \\ f \in \mathcal{Y}}} p(\boldsymbol{y}|\boldsymbol{x}) \log p(\boldsymbol{y}|\boldsymbol{x})$$

$$\sum_{\substack{\in \mathcal{Y} \\ f = 1}}^{T_y} \prod_{t=1}^{T_y} p(y_t|\boldsymbol{x})) (\sum_{t=1}^{T_y} \log p(y_t|\boldsymbol{x}))$$
Align table obtained from the alignment model
$$\sum_{\substack{i=1 \\ f \in \mathcal{Y} \\ f = 1}}^{T_y} \sum_{y_t \in \mathcal{A}(\boldsymbol{x})} p(y_t|\operatorname{Align}(y_t)) \log p(y_t|\operatorname{Align}(y_t))$$

$$= \frac{1}{|\mathcal{V}_x|} \sum_{x \in \mathcal{V}_x} \mathcal{H}(y|x)$$









Complexity (C(d)): 3.67

In practice, only measuring the complexity of the dataset is not enough for distillation data.

For distilled dataset, we also propose to measure the "faithfulness" which reflects to which extend, the distilled data is representative to the original parallel dataset. • We compute the **KL-divergence** of the alignment models between the real (r) and the

distilled dataset (d)

F(d) =





(b) NAT baseline.

(c) NAT trained on reduced mode by random selection. (d) NAT trained on distilled data set.

Complexity (C(d)): 3.30

Complexity (C(d)): 2.64

$$= \frac{1}{|\mathcal{V}_x|} \sum_{x \in \mathcal{V}_x} \sum_{y \in \mathcal{V}_y} p_r(y|x) \log \frac{p_r(y|x)}{p_d(y|x)}$$



26

We perform an extensive study over a variety of NAT and AT models with the proposed tools to analyze the **complexity** and **faithfulness** of the distilled dataset.

- Dataset: WMT14 English-German (En-De)
- Models and baseline scores (w/o distillation):



facebook Artificial Intelligence Research

```
German (En-De)
s (w/o distillation):
```

Models	Params	BLEU	Pass	Iters
AT models				
AT-tiny	16M	23.3	_	n
AT-small	37M	25.6	—	n
AT-base	65M	27.1	—	n
AT-big	218M	28.2		n
NAT models				
vanilla	71M	11.4	1	1
FlowSeq	73M	18.6	13	1
iNAT	66M	19.3	1	$k \ll n$
InsT	66M	20.9	1	$pprox \log_2 n$
MaskT	66M	23.5	1	10^{-1}
LevT	66M	25.2	1	$3k \ll n$
LevT-big	220M	26.5	≈ 3	$3k \ll n$



Analysis of the distilled dataset





• We visualize the complexity and faithfulness of our all 4 AT models (tiny, small, base, big) as well as the real data.

• As additional supporting metrics, we also plot the BLEU score (compared to the real data), showing it also correlates the data quality well.



Analysis of the distilled dataset



• As additional supporting metrics, we also plot the fuzzing reordering score for each dataset (Talbolt et al. 2011). A larger fuzzy reordering score indicates the more monotonic alignments.

> The distilled data looks much more monotonic to the English word order!

For more than 30 years, Josef Winkler has been writing from the heart, telling of the hardships of his childhood and youth. Seit mehr als 30 Jahren schreibt Josef Winkler aus dem Herzen und erzählt von der Not seiner Kindheit und Jugend . Josef Winkler schreibt sich seit mehr als 30 Jahren die Nöte seiner Kindheit und Jugend von der Seele .



Analysis of the distillation strategies

quality of distillation?





• In default, we take the beam-search output from the teacher model to create the distilled dataset. Will different decoding approaches affect the

YES. We must use beam-search (or at least greedy decoding).

Method	C(d)	F(d)	BLEU
	3.623	3.354	6.6
Тор 10)	2.411	2.932	14.6
	1.960	2.959	18.9
h	1.902	2.948	19.5



Analysis of the NAT models



weak

• Next, we show more results with different NAT models v.s. AT teachers are shown below. We always put the AT teacher scores (in red) for reference.



Analysis of the NAT models



weak

• The stronger the NAT model is, the closer it is to the AT teacher; • The teacher model does not have to be the upper-bound of the student (we will also come to this question later)



→ strong



Artificial Intelligence Research

weak

33



Improvements for WEAK student models

Born-Again Networks (BAN):



AT BLEU scores

Reborn Iterations

25

• Take the vanilla NAT model as an example.

Based on previous discussion, weak models require to be trained on simpler data. However, decreasing the size of the teacher model (e.g. base -> small) will hurt the faithfulness of the distilled data;

> distill the teacher model by its al output to train the







Improvements for STRONG student models

Sequence-level Interpolation (Seq-Inter):

- BLEU score from the ground-truth.

dbas bas

• Take the Levenshtein Transformer model as an example.

Based on previous discussion, strong models can be trained on more difficult data with high faithfulness. However, it requires training much stronger autoregressive teacher models (which is not easy); Kim & Rush, 2016 in fact also proposed improved version of distillation

named sequence-level interpolation, where we choose the K-best beam

search results and re-rank to select the sentences with the highest sentence-

	C(d)	F(d)	BLEU	
e inter	1.902	2.948	26.94	
	1.900	2.910	21.52	However, in practice this
				to the beam-size.



Implementation

Code for most of the NAT models can be found in Fairseq-py https://github.com/pytorch/fairseq/tree/master/examples/nonautoregressive_translation

facebook Artificial Intelligence Research

Revisiting Self-Training for Neural Sequence Generation

w/Junxian He, Jiajun Shen and Marc'Aurelio Ranzato

facebook Artificial Intelligence Research

ICLR2020









Self-Training



facebook Artificial Intelligence Research

• To answer the second question, we analyze how distillation works

when introducing more data. We keep teacher and student the same architecture.

In literature, such special setting of knowledge distillation is also called "self-training".

Different from the previous part, we usually need to "fine-tune" the student model on the real data (D) again (green arrow).

Furthermore, the fine-tuned student model can be treated as a new teacher, and we can repeat this loop multiple times, resulting in Iterative Self-Training.



Self-Training

How does self-training works in practice? • Test set BLEU on a subset of 100K parallel sentences from WMT14

English-German (En-De).

	Baseline	Iteration-1	Iteration-2	Iteration-3
Pseudo-train*	-	16.5	18.2	18.7
Train/Fine-tune	15.6	17.9	18.6	18.7



- The student trained only with distillation data, can usually outperform • its teacher!
- Fine-tuning on real data further boosts the translation quality, providing a better teacher model for the next iteration.



Even with the equal size teacher/student, the performance of the



The Secrets Behind Self-Training

We examine two possible hypotheses:

Decoding Strategy



- Typically, we always use "beam-search" instead of "sampling" from the teacher model's own distribution.
- The beam-searched targets serve as a "stronger" teacher model than • the student.

	Baseline	Beam-search	Sampling
Pseudo-train	-	16.5	16.1
Train/Fine-tune	15.6	17.9	17.0

only secrets behind the improvement.

The first possibility is that the gain comes from the "better" target.

The decoding strategy do affect the performance, however, is not the

40

The Secrets Behind Self-Training

We examine two possible hypotheses:

- **Decoding Strategy**
- **Noise during Training (Dropout)** •



The second assumption comes from the mismatched behaviors of "training" and "inference":

• during training \rightarrow self-training is not really "self".

	Baseline	Beam-search w/o Dropout	Sampling w/o Dropout	Beam-search	Sampling
Pseudo-train	_	15.8	15.5	16.5	16.1
Train/Fine-tune	15.6	16.3	16.0	17.9	17.0

Improvements disappeared on the pseudo-training phase!!

Dropouts are usually turned-off in the inference time, while open

41

Noisy Self-Training







Since we found "noise" during training useful, what if we add more? • Injecting synthetic noise in the input words, e.g. word swap, word deletion and word blanking (Lample et al., 2018).

	Baseline	Beam-search	Noisy Input + Beam-search
udo-train	_	16.5	16.6
in/Fine-tune	15.6	17.9	19.3

Injecting noise will not improve the pseudo-train results (should be expected as neither the source or the target are "REAL" sentences. However, injecting noise largely improve the performance on finetuning!

42

Noisy Self-Training



Since we found "noise" during training useful, what if we add more?

What is the **role** of "**noise**" in Self-training?

43

Case study on Toy Data



facebook Artificial Intelligence Research

it is always a good idea to look at some toy cases.

- Summing two integers in 0~99 as a sequence generation task; •
- Model works in the character level.
- We use only 250 pairs to training this task.

One good feature of this summing task is that we can easily visualize the results in a 2D space. For example:



When things are unclear and too difficult to explain in sequence generation,



44

Case study on Toy Data

Quantitative Analysis for Nois

Methods	smoothness	symmetric	error
baseline	9.1	9.8	7.6
ST	8.2	9.0	6.2
noisy ST	7.3	8.2	4.5

Table 3: Results on the toy sum dataset. For ST and noisy ST, smoothness (\downarrow) and symmetric (\downarrow) results are from the pseudo-training step, while test errors (\downarrow) are from fine-tuning, all at the first iteration.



Qualitative Analysis for Noisy Self-training



SV	Se	f _1	ra	ini	ina	ד*
Jy	JU				311	5

The injected noise will smooth the output space!



*Detailed definition of these metrics can be found in the paper.



45



We validate the proposed noisy self-training methods on both machine translation (MT) and text summarization (TS) tasks.

Machine Translation task:

- WMT14 English-German (En-De):
- FloRes English-Nepali (En-Ne)
- for comparison with target side monolingual data.

Mathada	WMT Engli	sh-German		FloRes English Nepali				
Methods	100K (+3.8M mono)	3.9M (+20M mono)	En-C	Drigin	Ne-Orig	in	Overall	
baseline	15.6	28.3	6	.7	2.3		4.8	
BT	20.5	_	8	.2	4.5		6.5	
noisy ST	21.4	29.3	8	.9	3.5		6.5	

```
simulated low-resource MT (100K) + 3.8M English (from the remaining)
full parallel data (3.9M) + 20M English (sampled from News Crawl)
```

```
real low-resource MT (560K) + 5M English (sampled from Wikipedia)
```

All noisy ST are performed 3 iterations. We also build up back-translation baselines

46



We validate the proposed noisy self-training methods on both machine translation (MT) and text summarization (TS) tasks.

Text Summarization:

- English Gigaword dataset simulated low-resource TS (100K, 640K); dataset)
- for comparison with target side summarizations.

Methods

MASS (Song et al., 2019)*

baseline BT noisy ST

full data (3.8M) + 4M monolingual documents (from the filtered Gigaword

All noisy ST are performed 3 iterations. We also build up back-translation baselines

100K	(+3.7M	(mono)	640K	(+3.2M	[mono)	3.8M (+4M mono)			
R 1	R2	RL	R 1	R2	RL	R 1	R2	RL	
_	_	_	_	_	_	38.7	19.7	36.0	
30.4	12.4	27.8	35.8	17.0	33.2	37.9	19.0	35.2	
32.2	13.8	29.6	37.3	18.4	34.6	_	_	_	
34.1	15.6	31.4	36.6	18.2	33.9	38.6	19.5	35.9	



Analysis of Dataset Size for Noisy Self-Training





Take the simulated WMT14 En-De data as an example:





Analysis of Noise-level injected in Noisy Self-Training

We vary the ratio of "word blanking" when injecting the noise.

Not surprisingly, the performance of selftraining drops a lot if the noise is too large.

• Take the simulated WMT14 En-De data as an example:







WAIT... one step back? What if we do not have new data, but inject noise onto parallel data?

facebook Artificial Intelligence Research Take the simulated WMT14 En-De data as an example:

• If following the same process as noisy self-training, only with parallel data still improves the performance (not as much as with However, if we only inject noise onto the source side, with real

sentence as the targets. The model will get much worse

ds	PT	FT
el baseline	-	15.6
ST, 100K mono + fake target	10.2	16.6
ST 3 8M mono + fake target	16.6	19 3
ST, 100K parallel + real target	6.7	11.3
ST, 100K parallel + fake target	10.4	16.0



Future works

Can we combine these two work?

- For instance, training a teacher AT model on limited parallel data;
- Distilled the model on much more monolingual data to train an NAT model

How can we get rid of distillation?

• For instance, GAN-style training for NAT models to handle multimodality

What is the best way to find the noise level for self-training?

• For instance, can we use meta-learning to learn to inject noise?

Some unrelated recent works...

Multilingual Denoising Pre-training for Neural Machine Translation

Yinhan Liu*, Jiatao Gu*, Naman Goyal*, Xian Li, Sergey Edunov Marjan Ghazvininejad, Mike Lewis, Luke Zettlemoyer Facebook AI Research

Facebook AI Research
{yinhanliu,jgu,naman,xianl,edunov
ghazvini,mikelewis,lsz}@fb.com



facebook Artificial Intelligence Research



Languages Data Source Size	En- WM 1(-Gu IT19)K	En- WM 91	- Kk IT19 I K	En-Vi IWSLT15 133K		En-Tr WMT17 207K		En-Ja IWSLT17 223K		En-Ko IWSLT17 230K	
Direction	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow
Random mBART25	0.0 0.3	0.0 0.1	0.8 7.4	0.2 2.5	23.6 36.1	24.8 35.4	12.2 22.5	9.5 17.8	10.4 19.1	12.3 19.4	15.3 24.6	16.3 22.6
Languages Data Source Size	En-Nl IWSLT17		En IWS 25	En-Ar IWSLT17		En-It IWSLT17		En-My WAT19 250K		-Ne oRes 4K	En-Ro WMT16	
Direction	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow
Random mBART25	34.6 43.3	29.3 34.8	27.5 37.6	16.9 21.6	31.7 39.8	28.0 34.0	23.3 28.3	34.9 36.9	7.6 14.5	4.3 7 4	34.0	34.3 37
Languages Data Source Size	En FLo	-Si Res	En IT	-Hi TB 6M	En WM	-Et IT18 4M	En WM	-Lt IT19				
Direction	\leftarrow	\rightarrow	+1.5 ←	\rightarrow	۰.1 ب	\rightarrow	\leftarrow	\rightarrow				bo
Random mBART25	7.2 13.7	1.2 3.3	10.9 23.5	14.2 20.8	22.6 27.8	17.9 21.4	18.1 22.4	17				ПE

Pre-traini	ng	Fine-tuning					
Model	Data	En→Ro	$Ro{\rightarrow}En$	+BT			
Random	None	34.3	34.0	36.8			
XLM (2019)	En Ro	-	35.6	38.5			
MASS (2019)	En Ro	-	-	39.1			
BART (2019)	En	-	-	38.0			
XLM-R (2019)	CC100	35.6	35.8	-			
BART-En	En	36.0	35.8	37.4			
BART-Ro	Ro	37.6	36.8	38.1			
mBART02	En Ro	38.5	38.5	39.9			
mBART25	CC25	37.7	37.8	38.8			

pa	per
	mo

	Model	Similar Pairs En-De En-Ro				Di En-	ar Pairs En-Si		
		\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow
facebook	Random XLM (2019)	21.0 34 3	17.2 26.4	19.4 31.8	21.2	0.0	0.0	0.0	0.0
Artificial Intelligence Research	MASS (2019)	35.2	28.3	33.1	35.2	-	-	-	-
	mBART	34.0	29.8	30.5	35.0	10.0	4.4	8.2	3.9



eck out our tomorrow

10⁶ (# of sentence pairs) 10^{7}

				VV			ng	Langua	iges				
							l ED	It TED	Ar TED	Hi News	Ne Wiki	Si Wiki	Gu Wiki
							7.0	6.8	6.2	7.2	4.2	5.9	0.0
		2				4.8	6.4	5.1	5.6	4.7	4.2	6.5	0.0
						8.5	9.5	9.1	8.7	9.6	8.8	11.1	0.0
						19.5	17.0	16.7	16.9	13.2	15.1	16.4	0.0
					<i>.</i> 0	37.8	22.3	21.6	22.6	16.4	18.5	22.1	0.0
			30.		21.2	27.0	43.3	34.1	31.0	24.6	23.3	27.3	0.0
			25.8	27.8	17.1	23.4	30.2	39.8	30.6	20.1	18.5	23.2	0.0
		ه	15.5	12.8	12.7	12.0	14.7	14.7	37.6	11.6	13.0	16.7	0.0
		3.2	10.1	9.9	5.8	6.7	6.1	5.0	7.6	23.5	14.5	13.0	0.0
L	Ne	2.1	6.7	6.5	5.0	4.3	3.0	2.2	5.2	17.9	14.5	10.8	0.0
	Si	5.0	5.7	3.8	3.8	1.3	0.9	0.5	3.5	8.1	8.9	13.7	0.0
	Gu	8.2	8.5	4.7	5.4	3.5	2.1	0.0	6.2	13.8	13.5	12.8	0.3



facebook Artificial Intelligence Research

Thank you!

We are also hiring Research Interns, AI Residents and Full-time Researchers at FAIR! Let me know if you are interested! jgu@fb.com