

## 89688: Statistical Machine Translation

#### **Unsupervised Neural Machine Translation**

Roee Aharoni Computer Science Department Bar Ilan University

June 2020

 NMT is better than SMT only when given >10m parallel words

**BLEU Scores with Varying Amounts of Training Data** 

30  $20 \begin{array}{c} 26.2 \\ 24.9 \\ 23.4 \\ 23.4 \\ 21.2 \\ 22.2 \\ 23.5 \\ 21.2 \\ 22.2 \\ 23.5 \\ 21.2 \\ 22.2 \\ 23.5 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 23.5 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 21.2 \\ 22.2 \\ 24.9 \\ 24.9 \\ 21.2 \\ 24.9 \\ 24.9 \\ 21.2 \\ 24.9 \\$ 10 — Phrase-Based with Big LM Phrase-Based -----Neural -0- $10^{6}$  $10^{8}$ 10 Corpus Size (English Words)

- NMT is better than SMT only when given >10m parallel words
- NMT is better than "Semi Supervised" SMT (SMT + a large language model) only when given >100m parallel words

**BLEU Scores with Varying Amounts of Training Data** 



- NMT is better than SMT only when given >10m parallel words
- NMT is better than "Semi Supervised" SMT (SMT + a large language model) only when given >100m parallel words
- But getting parallel data is expensive!

**BLEU Scores with Varying Amounts of Training Data** 



- NMT is better than SMT only when given >10m parallel words
- NMT is better than "Semi Supervised" SMT (SMT + a large language model) only when given >100m parallel words
- But getting parallel data is expensive!
- Can we do well using only **monolingual data?**

**BLEU Scores with Varying Amounts of Training Data** 



• "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013

- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation



- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success







- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success
- Weakly supervised requires a small dictionary (5000 entries)







• Both submitted to ICLR 2018 with critical acclaim (October 2017)

- Both submitted to ICLR 2018 with critical acclaim (October 2017)
- Similar motivations both try to tackle:

- Both submitted to ICLR 2018 with critical acclaim (October 2017)
- Similar motivations both try to tackle:
  - Structure/Fluency how to determine the correct word order in the output?

- Both submitted to ICLR 2018 with critical acclaim (October 2017)
- Similar motivations both try to tackle:
  - Structure/Fluency how to determine the correct word order in the output?
  - Semantics/Adequacy how to pick the correct translations given the source?

- Both submitted to ICLR 2018 with critical acclaim (October 2017)
- Similar motivations both try to tackle:
  - Structure/Fluency how to determine the correct word order in the output?
  - Semantics/Adequacy how to pick the correct translations given the source?
- Very similar modeling tricks (with slight differences)



- Model Architecture:
  - Shared GRU encoder, Separate GRU decoders



- Shared GRU encoder, Separate GRU decoders
- Attention



- Shared GRU encoder, Separate GRU decoders
- Attention
- Main "Tricks":



- Shared GRU encoder, Separate GRU decoders
- Attention
- Main "Tricks":
  - Fixed, unsupervised cross-lingual embeddings (Adequacy)



- Shared GRU encoder, Separate GRU decoders
- Attention
- Main "Tricks":
  - Fixed, unsupervised cross-lingual embeddings (Adequacy)
  - Backtranslation loss (**Adequacy**)



- Shared GRU encoder, Separate GRU decoders
- Attention
- Main "Tricks":
  - Fixed, unsupervised cross-lingual embeddings (Adequacy)
  - Backtranslation loss (Adequacy)
  - Denoising auto-encoder loss (Fluency)



• Artetxe, Labake & Agirre, ACL 2017

- Artetxe, Labake & Agirre, ACL 2017
- Start with monolingual embedding spaces in two languages (trained using word2vec)



From Artetxe, ACL 2017



- Artetxe, Labake & Agirre, ACL 2017
- Start with monolingual embedding spaces in two languages (trained using word2vec)
- Learn a linear mapping from one language to the other:



From Artetxe, ACL 2017

- Artetxe, Labake & Agirre, ACL 2017
- Start with monolingual embedding spaces in two languages (trained using word2vec)
- Learn a linear mapping from one language to the other:
  - Start with a **seed dictionary**. Clever idea: use numerals (5-5, 1989-1989...) as seed dictionary fully unsupervised
  - Optimize the mapping W w.r.t the dictionary:  $\arg \min_{W \in O(n)} \sum_{i} ||X_{i*}W Z_{j*}||^2$



From Artetxe, ACL 2017



25

- Artetxe, Labake & Agirre, ACL 2017
- Start with monolingual embedding spaces in two languages (trained using word2vec)
- Learn a linear mapping from one language to the other:
  - Start with a **seed dictionary**. Clever idea: use numerals (5-5, 1989-1989...) as seed dictionary fully unsupervised
  - Optimize the mapping W w.r.t the dictionary:  $\arg \min_{W \in O(n)} \sum_{i} ||X_{i*}W Z_{j*}||^2$
- Extract a new dictionary and **repeat iteratively** until a convergence threshold is met



From Artetxe, ACL 2017



25

- Artetxe, Labake & Agirre, ACL 2017
- Start with monolingual embedding spaces in two languages (trained using word2vec)
- Learn a linear mapping from one language to the other:
  - Start with a **seed dictionary**. Clever idea: use numerals (5-5, 1989-1989...) as seed dictionary fully unsupervised
  - Optimize the mapping W w.r.t the dictionary:  $\arg \min_{W \in O(n)} \sum_{i} ||X_{i*}W Z_{j*}||^2$
- Extract a new dictionary and **repeat iteratively** until a convergence threshold is met



From Artetxe, ACL 2017



25

- Artetxe, Labake & Agirre, ACL 2017
- Start with monolingual embedding spaces in two languages (trained using word2vec)
- Learn a linear mapping from one language to the other:
  - Start with a **seed dictionary**. Clever idea: use numerals (5-5, 1989-1989...) as seed dictionary - fully unsupervised
  - Optimize the mapping W w.r.t the dictionary:  $\arg \min_{W \in O(n)} \sum_{i} ||X_{i*}W Z_{j*}||^2$
- Extract a new dictionary and **repeat iteratively** until a convergence threshold is met



From Artetxe, ACL 2017

- Artetxe, Labake & Agirre, ACL 2017
- Start with monolingual embedding spaces in two languages (trained using word2vec)
- Learn a linear mapping from one language to the other:
  - Start with a **seed dictionary**. Clever idea: use numerals (5-5, 1989-1989...) as seed dictionary - fully unsupervised
  - Optimize the mapping W w.r.t the dictionary:  $\arg \min_{W \in O(n)} \sum_{i} ||X_{i*}W Z_{j*}||^2$
- Extract a new dictionary and **repeat iteratively** until a convergence threshold is met



From Artetxe, ACL 2017

#### Learning Semantics: Back-Translation

### Learning Semantics: Back-Translation

- We need to learn a mapping from one language (L1) into another (L2), but we don't have parallel data
- We need to learn a mapping from one language (L1) into another (L2), but we don't have parallel data
- Solution: create **synthetic** parallel data by translating with the current model (possible since the model is bidirectional)

- We need to learn a mapping from one language (L1) into another (L2), but we don't have parallel data
- Solution: create **synthetic** parallel data by translating with the current model (possible since the model is bidirectional)
- Use the synthetic data for training using cross entropy loss

- We need to learn a mapping from one language (L1) into another (L2), but we don't have parallel data
- Solution: create **synthetic** parallel data by translating with the current model (possible since the model is bidirectional)
- Use the synthetic data for training using cross entropy loss

L1 Sentence

- We need to learn a mapping from one language (L1) into another (L2), but we don't have parallel data
- Solution: create **synthetic** parallel data by translating with the current model (possible since the model is bidirectional)
- Use the synthetic data for training using cross entropy loss



translate using the current model





- We need to learn a mapping from one language (L1) into another (L2), but we don't have parallel data
- Solution: create **synthetic** parallel data by translating with the current model (possible since the model is bidirectional)
- Use the synthetic data for training using cross entropy loss



- We need to learn a mapping from one language (L1) into another (L2), but we don't have parallel data
- Solution: create **synthetic** parallel data by translating with the current model (possible since the model is bidirectional)
- Use the synthetic data for training using cross entropy loss
- This is not entirely useless since the cross-lingual embeddings do carry some alignment signal



• The decoder needs to learn how to organize the words on the target side

- The decoder needs to learn how to organize the words on the target side
- We could train it to predict a sentence given itself auto-encoding

- The decoder needs to learn how to organize the words on the target side
- We could train it to predict a sentence given itself auto-encoding
  - But this would lead it to learn trivial copying!

- The decoder needs to learn how to organize the words on the target side
- We could train it to predict a sentence given itself auto-encoding
  - But this would lead it to learn trivial copying!
- Introduce "noise" by randomly swapping adjacent words (N/2 times) in the input, to force the decoder to learn word ordering

- The decoder needs to learn how to organize the words on the target side
- We could train it to predict a sentence given itself auto-encoding
  - But this would lead it to learn trivial copying!
- Introduce "noise" by randomly swapping adjacent words (N/2 times) in the input, to force the decoder to learn word ordering

cat The on sat mat the

The cat sat on the mat

- The decoder needs to learn how to organize the words on the target side
- We could train it to predict a sentence given itself auto-encoding
  - But this would lead it to learn trivial copying!
- Introduce "noise" by randomly swapping adjacent words (N/2 times) in the input, to force the decoder to learn word ordering

cat The on sat mat the

The cat sat on the mat

manger J'aime croissants des

- The decoder needs to learn how to organize the words on the target side
- We could train it to predict a sentence given itself auto-encoding
  - But this would lead it to learn trivial copying!
- Introduce "noise" by randomly swapping adjacent words (N/2 times) in the input, to force the decoder to learn word ordering
- Train using conventional cross entropy loss

cat The on sat mat the

The cat sat on the mat

manger J'aime croissants des

 Training iterations alternate between denoising and back-translation



- Training iterations alternate between denoising and back-translation
- Each training iteration is composed of:



- Training iterations alternate between denoising and back-translation
- Each training iteration is composed of:
  - Denoising batch: L1 to L1



- Training iterations alternate between denoising and back-translation
- Each training iteration is composed of:
  - Denoising batch: L1 to L1
  - Denoising batch: L2 to L2



- Training iterations alternate between denoising and back-translation
- Each training iteration is composed of:
  - Denoising batch: L1 to L1
  - Denoising batch: L2 to L2
  - Back-translation batch: L1 to L2



- Training iterations alternate between denoising and back-translation
- Each training iteration is composed of:
  - Denoising batch: L1 to L1
  - Denoising batch: L2 to L2
  - Back-translation batch: L1 to L2
  - Back-translation batch: L2 to L1



- Training iterations alternate between denoising and back-translation
- Each training iteration is composed of:
  - Denoising batch: L1 to L1
  - Denoising batch: L2 to L2
  - Back-translation batch: L1 to L2
  - Back-translation batch: L2 to L1
- When do we stop? Can't use parallel validation set!



- Training iterations alternate between denoising and back-translation
- Each training iteration is composed of:
  - Denoising batch: L1 to L1
  - Denoising batch: L2 to L2
  - Back-translation batch: L1 to L2
  - Back-translation batch: L2 to L1
- When do we stop? Can't use parallel validation set!
  - Train for a fixed amount of iterations (300k)



 Denoising alone is weaker than the nearestneighbor baseline

		FR-EN	EN-FR	DE-EN	EN
Unsupervised	<ol> <li>Baseline (emb. nearest neighbor)</li> <li>Proposed (denoising)</li> <li>Proposed (+ backtranslation)</li> <li>Proposed (+ BPE)</li> </ol>	9.98 7.28 15.56 15.56	6.25 5.33 15.13 14.36	7.07 3.64 10.21 10.16	4 2 6 6
Semi-supervised	5. Proposed (full) + 100k parallel	21.81	21.74	15.24	10
Supervised	<ol> <li>Comparable NMT</li> <li>GNMT (Wu et al., 2016)</li> </ol>	20.48	19.89 38.95	15.04 -	11 24



- Denoising alone is weaker than the nearestneighbor baseline
- Denoising+Back-translation significantly improves results

		FR-EN	EN-FR	DE-EN	EN
Unsupervised	<ol> <li>Baseline (emb. nearest neighbor)</li> <li>Proposed (denoising)</li> <li>Proposed (+ backtranslation)</li> <li>Proposed (+ BPE)</li> </ol>	9.98 7.28 15.56 15.56	6.25 5.33 15.13 14.36	7.07 3.64 10.21 10.16	4 2 6 6
Semi-supervised	5. Proposed (full) + 100k parallel	21.81	21.74	15.24	10
Supervised	<ol> <li>Comparable NMT</li> <li>GNMT (Wu et al., 2016)</li> </ol>	20.48	19.89 38.95	15.04 -	11 24



- Denoising alone is weaker than the nearestneighbor baseline
- Denoising+Back-translation significantly improves results
- No clear benefit from BPE (harder to learn embeddings?)

		FR-EN	EN-FR	DE-EN	EN
Unsupervised	<ol> <li>Baseline (emb. nearest neighbor)</li> <li>Proposed (denoising)</li> <li>Proposed (+ backtranslation)</li> <li>Proposed (+ BPE)</li> </ol>	9.98 7.28 15.56 15.56	6.25 5.33 15.13 14.36	7.07 3.64 10.21 10.16	4 2 6 6
Semi-supervised	5. Proposed (full) + 100k parallel	21.81	21.74	15.24	10
Supervised	<ol> <li>Comparable NMT</li> <li>GNMT (Wu et al., 2016)</li> </ol>	20.48	19.89 38.95	15.04 -	11 24



- Denoising alone is weaker than the nearestneighbor baseline
- Denoising+Back-translation significantly improves results
- No clear benefit from BPE (harder to learn embeddings?)
- Semi supervised learning can also use this framework with notable gains
- Still a very large gap from the supervised approach (but a nice start nonetheless)

		FR-EN	EN-FR	DE-EN	EN
Unsupervised	<ol> <li>Baseline (emb. nearest neighbor)</li> <li>Proposed (denoising)</li> <li>Proposed (+ backtranslation)</li> <li>Proposed (+ BPE)</li> </ol>	9.98 7.28 15.56 15.56	6.25 5.33 15.13 14.36	7.07 3.64 10.21 10.16	4 2 6 6
Semi-supervised	5. Proposed (full) + 100k parallel	21.81	21.74	15.24	10
Supervised	6. Comparable NMT 7. GNMT (Wu et al., 2016)	20.48	19.89 38.95	15.04	11 24







• Model Architecture:

• Shared GRU encoder, Shared GRU decoder





- Shared GRU encoder, Shared GRU decoder
- Attention





- Shared GRU encoder, Shared GRU decoder
- Attention
- Main "Tricks":





- Shared GRU encoder, Shared GRU decoder
- Attention
- Main "Tricks":
  - Changing, adversarially trained unsupervised cross-lingual embeddings (**Adequacy**)





- Shared GRU encoder, Shared GRU decoder
- Attention
- Main "Tricks":
  - Changing, adversarially trained unsupervised cross-lingual embeddings (**Adequacy**)
  - Backtranslation loss (**Adequacy**)





- Shared GRU encoder, Shared GRU decoder
- Attention
- Main "Tricks":
  - Changing, adversarially trained unsupervised cross-lingual embeddings (Adequacy)
  - Backtranslation loss (**Adequacy**)
  - Denoising auto-encoder loss (**Fluency**)




#### Paper II: Lample, Denoyer and Ranzato (FAIR)

#### • Model Architecture:

- Shared GRU encoder, Shared GRU decoder
- Attention
- Main "Tricks":
  - Changing, adversarially trained unsupervised cross-lingual embeddings (Adequacy)
  - Backtranslation loss (**Adequacy**)
  - Denoising auto-encoder loss (**Fluency**)
  - Adversarial loss





• Introduced by Ganin et al., 2016 for domain adaption in computer vision

- Introduced by Ganin et al., 2016 for domain adaption in computer vision
- The general idea: force the model to "unlearn" a specific objective to make it learn better representation for the target objective

- Introduced by Ganin et al., 2016 for domain adaption in computer vision
- The general idea: force the model to "unlearn" a specific objective to make it learn better representation for the target objective
- Used twice here:

- Introduced by Ganin et al., 2016 for domain adaption in computer vision
- The general idea: force the model to "unlearn" a specific objective to make it learn better representation for the target objective
- Used twice here:
  - In the cross-lingual embedding learning to learn a mapping from one embedding space to the other:



Conneau et al. 2017

- Introduced by Ganin et al., 2016 for domain adaption in computer vision
- The general idea: force the model to "unlearn" a specific objective to make it learn better representation for the target objective
- Used twice here:
  - In the cross-lingual embedding learning to learn a mapping from one embedding space to the other:
  - In the NMT training to "push" the representations from the two languages to a shared "semantic" space





Conneau et al. 2017

$$p_D(l|z_1,...,z_m) \propto \prod_{j=1}^m p_D(\ell|z_j),$$

 $\mathcal{L}_{adv}(\theta_{enc}, \mathcal{Z}|\theta_D) = -\mathbb{E}_{(x_i, \ell_i)}[\log p_D(\ell_j|e(x_i, \ell_i))]$ 



 When do we stop training without a validation set? can we do better than fixed amount of updates?

- When do we stop training without a validation set? can we do better than fixed amount of updates?
- Measure "corruption" when translating a sentence back and forth using the model (in both directions), using BLEU

 $MS(e, d, \mathcal{D}_{src}, \mathcal{D}_{tgt}) =$ 

 $\frac{1}{2} \mathbb{E}_{x \sim \mathcal{D}_{src}} \left[ BLEU(x, M_{src \to tgt} \circ M_{tgt \to src}(x)) \right] + \\ \frac{1}{2} \mathbb{E}_{x \sim \mathcal{D}_{tgt}} \left[ BLEU(x, M_{tgt \to src} \circ M_{src \to tgt}(x)) \right]$ 

- When do we stop training without a validation set? can we do better than fixed amount of updates?
- Measure "corruption" when translating a sentence back and forth using the model (in both directions), using BLEU
- Correlates well with "supervised" BLEU, no need for parallel sentences

 $MS(e,d,\mathcal{D}_{src},\mathcal{D}_{tgt})$  =

 $\frac{1}{2} \mathbb{E}_{x \sim \mathcal{D}_{src}} \left[ BLEU(x, M_{src \to tgt} \circ M_{tgt \to src}(x)) \right] + \frac{1}{2} \mathbb{E}_{x \sim \mathcal{D}_{tgt}} \left[ BLEU(x, M_{tgt \to src} \circ M_{src \to tgt}(x)) \right]$ 



 Model significantly outperforms word-byword baselines, showing the importance of the back-translation + denoising + adversarial approach

	Multi30k-Task1					WMT	
	en-fr	fr-en	de-en	en-de	en-fr	fr-en	de-en
Supervised	56.83	50.77	38.38	35.16	27.97	26.13	25.61
word-by-word word reordering oracle word reordering	8.54 - 11.62	16.77 - 24.88	15.72 - 18.27	5.39 - 6.79	6.28 6.68 10.12	10.09 11.69 20.64	10.77 10.84 19.42
Our model: 1st iteration Our model: 2nd iteration Our model: 3rd iteration	27.48 31.72 32.76	28.07 30.49 32.07	23.69 24.73 26.26	19.32 21.16 22.74	12.10 14.42 15.05	11.79 13.49 14.31	11.10 13.25 13.33



- Model significantly outperforms word-byword baselines, showing the importance of the back-translation + denoising + adversarial approach
- Supervised models are still significantly better

	Multi30k-Task1					WMT	
	en-fr	fr-en	de-en	en-de	en-fr	fr-en	de-en
Supervised	56.83	50.77	38.38	35.16	27.97	26.13	25.61
word-by-word word reordering oracle word reordering	8.54 - 11.62	16.77 - 24.88	15.72 - 18.27	5.39 - 6.79	6.28 6.68 10.12	10.09 11.69 20.64	10.77 10.84 19.42
Our model: 1st iteration Our model: 2nd iteration Our model: 3rd iteration	27.48 31.72 32.76	28.07 30.49 32.07	23.69 24.73 26.26	19.32 21.16 22.74	12.10 14.42 15.05	11.79 13.49 14.31	11.10 13.25 13.33



- Model significantly outperforms word-byword baselines, showing the importance of the back-translation + denoising + adversarial approach
- Supervised models are still significantly better
- Unsupervised models performance is equivalent to a supervised model with ~100k parallel sentences

	Multi30k-Task1					WMT	
	en-fr	fr-en	de-en	en-de	en-fr	fr-en	de-en
Supervised	56.83	50.77	38.38	35.16	27.97	26.13	25.61
word-by-word word reordering oracle word reordering	8.54 - 11.62	16.77 - 24.88	15.72 - 18.27	5.39 - 6.79	6.28 6.68 10.12	10.09 11.69 20.64	10.77 10.84 19.42
Our model: 1st iteration Our model: 2nd iteration Our model: 3rd iteration	27.48 31.72 32.76	28.07 30.49 32.07	23.69 24.73 26.26	19.32 21.16 22.74	12.10 14.42 15.05	11.79 13.49 14.31	11.10 13.25 13.33







Back-translation, pre-trained word vectors and de-noising are crucial

- Back-translation, pre-trained word vectors and de-noising are crucial
- Adversarial loss gives a nice boost of ~3-6 points

- Back-translation, pre-trained word vectors and de-noising are crucial
- Adversarial loss gives a nice boost of ~3-6 points
- Best model obtained using all components together

- Back-translation, pre-trained word vectors and de-noising are crucial
- Adversarial loss gives a nice boost o ~3-6 points
- Best model obtained using all components together

	en-fr	fr-en	de-en	en
$\lambda_{cd} = 0$	25.44	27.14	20.56	14
Without pretraining	25.29	26.10	21.44	17
Without pretraining, $\lambda_{cd} = 0$	8.78	9.15	7.52	6
Without noise, $C(x) = x$	16.76	16.85	16.85	14
$\lambda_{auto} = 0$	24.32	20.02	19.10	14
$\lambda_{adv} = 0$	24.12	22.74	19.87	15
Full	27.48	28.07	23.69	19



•Both models:

#### •Both models:

•Heavily rely on bilingual word embeddings

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

#### •Both models:

•Heavily rely on bilingual word embeddings

•Heavily rely on de-noising and back-translation using a shared encoder

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

#### •Notable Differences:

•Perform back-translation once per epoch (Lample et al.) vs. after every update (Artetxe et al.)

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

- •Perform back-translation once per epoch (Lample et al.) vs. after every update (Artetxe et al.)
- •BPE and word-based modeling (Artetxe et al.) vs. word based alone (Lample et al.)

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

- •Perform back-translation once per epoch (Lample et al.) vs. after every update (Artetxe et al.)
- •BPE and word-based modeling (Artetxe et al.) vs. word based alone (Lample et al.)
- Fixed embeddings (Artetxe et al.) vs. changing (Lample et al.)

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

- •Perform back-translation once per epoch (Lample et al.) vs. after every update (Artetxe et al.)
- •BPE and word-based modeling (Artetxe et al.) vs. word based alone (Lample et al.)
- Fixed embeddings (Artetxe et al.) vs. changing (Lample et al.)
- Different decoder per language (Artetxe et al.) vs. shared encoder (Lample et al.)

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

- •Perform back-translation once per epoch (Lample et al.) vs. after every update (Artetxe et al.)
- •BPE and word-based modeling (Artetxe et al.) vs. word based alone (Lample et al.)
- Fixed embeddings (Artetxe et al.) vs. changing (Lample et al.)
- Different decoder per language (Artetxe et al.) vs. shared encoder (Lample et al.)
- •Adversarial training (Lample et al.)

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

#### •Notable Differences:

- •Perform back-translation once per epoch (Lample et al.) vs. after every update (Artetxe et al.)
- •BPE and word-based modeling (Artetxe et al.) vs. word based alone (Lample et al.)
- Fixed embeddings (Artetxe et al.) vs. changing (Lample et al.)
- Different decoder per language (Artetxe et al.) vs. shared encoder (Lample et al.)
- •Adversarial training (Lample et al.)

•Slightly different noise method (Lample et al.) - swapping and dropping words, also adding noise before back-translation

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

#### •Notable Differences:

- •Perform back-translation once per epoch (Lample et al.) vs. after every update (Artetxe et al.)
- •BPE and word-based modeling (Artetxe et al.) vs. word based alone (Lample et al.)
- Fixed embeddings (Artetxe et al.) vs. changing (Lample et al.)
- Different decoder per language (Artetxe et al.) vs. shared encoder (Lample et al.)
- •Adversarial training (Lample et al.)
- •Use unsupervised model selection criterion (Lample et al.) vs. fixed amount of updates

•Slightly different noise method (Lample et al.) - swapping and dropping words, also adding noise before back-translation

#### •Both models:

- •Heavily rely on bilingual word embeddings
- •Heavily rely on de-noising and back-translation using a shared encoder

#### •Notable Differences:

- •Perform back-translation once per epoch (Lample et al.) vs. after every update (Artetxe et al.)
- •BPE and word-based modeling (Artetxe et al.) vs. word based alone (Lample et al.)
- Fixed embeddings (Artetxe et al.) vs. changing (Lample et al.)
- Different decoder per language (Artetxe et al.) vs. shared encoder (Lample et al.)
- •Adversarial training (Lample et al.)
- •Use unsupervised model selection criterion (Lample et al.) vs. fixed amount of updates
- •Initialize back-translation using nearest-neighbor word-by-word translation (Lample et al.)

•Slightly different noise method (Lample et al.) - swapping and dropping words, also adding noise before back-translation

### Phrase-Based Unsupervised NMT

# Phrase-Based Unsupervised NMT

 A second wave of works integrated phrasebased models for unsupervised NMT (also from the same authors):

### Phrase-Based Unsupervised NMT

- A second wave of works integrated phrasebased models for unsupervised NMT (also from the same authors):
  - Lample et al. (2018)

#### Phrase-Based & Neural Unsupervised Machine Translation

Guillaume Lample<sup>†</sup> Facebook AI Research Sorbonne Universités glample@fb.com Myle Ott Facebook AI Research myleott@fb.com

Alexis Conneau Facebook AI Research Université Le Mans aconneau@fb.com

Ludovic Denoyer<sup>†</sup> Sorbonne Universités ludovic.denoyer@lip6.fr Marc'Aurelio Ranzato Facebook AI Research ranzato@fb.com
## Phrase-Based Unsupervised NMT

- A second wave of works integrated phrasebased models for unsupervised NMT (also from the same authors):
  - Lample et al. (2018)
  - Artetxe et al. (2018)

#### **Phrase-Based & Neural Unsupervised Machine Translation**

Guillaume Lample<sup>†</sup> Facebook AI Research Sorbonne Universités glample@fb.com Myle Ott Facebook AI Research myleott@fb.com

Alexis Conneau Facebook AI Research Université Le Mans aconneau@fb.com

Ludovic Denoyer<sup>†</sup> Sorbonne Universités ludovic.denoyer@lip6.fr Marc'Aurelio Ranzato Facebook AI Research ranzato@fb.com

#### **Unsupervised Statistical Machine Translation**

## Phrase-Based Unsupervised NMT

- A second wave of works integrated phrasebased models for unsupervised NMT (also from the same authors):
  - Lample et al. (2018)
  - Artetxe et al. (2018)
  - Artetxe et al (2019)

#### **Phrase-Based & Neural Unsupervised Machine Translation**

Guillaume Lample<sup>†</sup> Facebook AI Research Sorbonne Universités glample@fb.com Myle Ott Facebook AI Research myleott@fb.com

Alexis Conneau Facebook AI Research Université Le Mans aconneau@fb.com

Ludovic Denoyer<sup>†</sup> Sorbonne Universités ludovic.denoyer@lip6.fr Marc'Aurelio Ranzato Facebook AI Research ranzato@fb.com

#### **Unsupervised Statistical Machine Translation**

Mikel Artetxe, Gorka Labaka, Eneko Agirre IXA NLP Group University of the Basque Country (UPV/EHU) {mikel.artetxe,gorka.labaka,e.agirre}@ehu.eus

#### An Effective Approach to Unsupervised Machine Translation

## Phrase-Based Unsupervised NMT

- A second wave of works integrated phrasebased models for unsupervised NMT (also from the same authors):
  - Lample et al. (2018)
  - Artetxe et al. (2018)
  - Artetxe et al (2019)
- Makes sense, as SMT was shown to work better in low resource scenarios

#### **Phrase-Based & Neural Unsupervised Machine Translation**

Guillaume Lample<sup>†</sup> Facebook AI Research Sorbonne Universités glample@fb.com Myle Ott Facebook AI Research myleott@fb.com

Alexis Conneau Facebook AI Research Université Le Mans aconneau@fb.com

Ludovic Denoyer<sup>†</sup> Sorbonne Universités ludovic.denoyer@lip6.fr Marc'Aurelio Ranzato Facebook AI Research ranzato@fb.com

#### **Unsupervised Statistical Machine Translation**

Mikel Artetxe, Gorka Labaka, Eneko Agirre IXA NLP Group University of the Basque Country (UPV/EHU) {mikel.artetxe,gorka.labaka,e.agirre}@ehu.eus

#### An Effective Approach to Unsupervised Machine Translation

• Artetxe, Labaka and Agirre, 2018

**Unsupervised Statistical Machine Translation** 



- Artetxe, Labaka and Agirre, 2018
- Main idea: use SMT instead of NMT

#### **Unsupervised Statistical Machine Translation**



- Artetxe, Labaka and Agirre, 2018
- Main idea: use SMT instead of NMT
- Train n-gram embeddings using a variation of skip-gram

#### **Unsupervised Statistical Machine Translation**



- Artetxe, Labaka and Agirre, 2018
- Main idea: use SMT instead of NMT
- Train n-gram embeddings using a variation of skip-gram
- Learn a mapping between the embedding spaces

#### **Unsupervised Statistical Machine Translation**



- Artetxe, Labaka and Agirre, 2018
- Main idea: use SMT instead of NMT
- Train n-gram embeddings using a variation of skip-gram
- Learn a mapping between the embedding spaces
- Create a phrase table by computing translation probabilities using softmax over the cosinesimilarities to the 100 nearest neighbours

#### **Unsupervised Statistical Machine Translation**



- Artetxe, Labaka and Agirre, 2018
- Main idea: use SMT instead of NMT
- Train n-gram embeddings using a variation of skip-gram
- Learn a mapping between the embedding spaces
- Create a phrase table by computing translation probabilities using softmax over the cosinesimilarities to the 100 nearest neighbours
- Tune the resulting system using iterative backtranslation

#### **Unsupervised Statistical Machine Translation**



 Much better than previous unsupervised NMT approaches across 6 language pairs

		WMT-16				
	FR-EN	EN-FR	DE-EN	EN-DE	DE-EN	EN
Artetxe et al. (2018c)	15.56	15.13	10.21	6.55	-	
Lample et al. (2018)	14.31	15.05	-	-	13.33	9.
Yang et al. (2018)	15.58	16.97	-	-	14.62	10
Proposed system	25.87	26.22	17.43	14.08	23.05	18



- Much better than previous unsupervised NMT approaches across 6 language pairs
- Unsupervised Tuning and Iterative Back-Translation are important

	WMT-14							WMT-16		
-	FR	-EN 1	EN-FR	DE	E-EN	EN-DE		DE-EN	EN	
Artetxe et al. (2018c)	15	.56	15.13	10	0.21	6.55		-		
Lample et al. (2018)	14	.31	15.05		-	-		13.33	9.	
Yang et al. (2018)	15	.58	16.97		-	-		14.62	10.	
Proposed system	25	.87	26.22	17	7.43	14.08		23.05	18	
	WMT-14							WMT-16		
	_	FR-EN	EN-F	R I	DE-EN	EN-DE	2	DE-EN	EN	
Unsupervised SMT		21.16	20.13	3	13.86	10.59		18.01	13	
+ unsupervised tuning		22.17	22.22	2	14.73	10.64		18.21	13	
+ iterative refinement (it1)		24.81	26.53	3	16.01	13.45		20.76	16	
+ iterative refinement (it2)		26.13	26.57	7	17.30	13.95		22.80	18	
+ iterative refinement (it3)		25.87	26.22	2	17.43	14.08		23.05	18	



- Much better than previous unsupervised NMT approaches across 6 language pairs
- Unsupervised Tuning and Iterative Back-Translation are important
- Still far from the supervised approaches

		WMT-14						WMT-16		
	-	FR-EN	EN	-FR	DE-	EN	EN-	DE	DE-EN	EN
Artetxe et a	l. (2018c)	15.56	15	5.13	10.2	21	6.5	55	-	-
Lample et a	1. (2018)	14.31	15	5.05	-		-		13.33	9.0
Yang et al.	(2018)	15.58	16	5.97	-	-			14.62	10.
Proposed sy	vstem	25.87	20	5.22	17.4	43	14.	08	23.05	18
				N	/MT-1	4			WM	<b>T-16</b>
		FR-	EN	EN-F	R D	E-EN	EN	N-DE	DE-EN	EN·
Unsupervise	ed SMT	21.	16	20.13	3 1	3.86	1	0.59	18.01	13
+ unsupervised tuning 22		22.	.17 22.		2 1	14.73		0.64	18.21	13.
+ iterative refinement (it1) 2		1) 24.	81	26.53	3 1	16.01		3.45	20.76	16.
+ iterative r	efinement (it	2) <b>26</b> .	13	26.57	/ 1	7.30	1.	3.95	22.80	18
+ iterative r	efinement (it	3) 25.	.87	26.22	2 1	7.43	14	4.08	23.05	18
			WMT-14					WI	MT-1	
			FR-I	EN E	N-FR	DE-	EN	EN-DE	DE-EN	EN
	NMT (transf	former)	-		41.8	-		28.4	-	
Supervised	WMT best		35.	0	35.8	29.	.0	20.6	40.2	3
	SMT (europ	arl)	30.6	51 3	30.82	20.	83	16.60	26.38	22
	+ w/o lexica	al reord. 30		54 3	30.33	20.	37	16.34	25.99	22
	+ constraine	d vocab.	30.0	)4 3	30.10	19.	9.91 16.32		25.66	2
	+ unsup. tur	ing	29.3	32 2	29.46 17.		75	15.45	23.35	19
Unsup.	Proposed sy	stem	25.8	25.87 26		17.43		14.08	23.05	18



 First proposed by Lample et al. (2018) proposed similar ideas with SMT, and a joint approach by training NMT on SMT outputs

#### **Phrase-Based & Neural Unsupervised Machine Translation**

Guillaume Lample<sup>†</sup> Facebook AI Research Sorbonne Universités glample@fb.com Myle Ott Facebook AI Research myleott@fb.com

Alexis Conneau Facebook AI Research Université Le Mans aconneau@fb.com

Ludovic Denoyer<sup>†</sup> Sorbonne Universités ludovic.denoyer@lip6.fr Marc'Aurelio Ranzato Facebook AI Research ranzato@fb.com



- First proposed by Lample et al. (2018) proposed similar ideas with SMT, and a joint approach by training NMT on SMT outputs
- Improved by Artetxe et al (2019) with better tuning of the SMT model and a gradual mixing of SMT and NMT backtranslations

#### **Phrase-Based & Neural Unsupervised Machine Translation**

Guillaume Lample<sup>†</sup> Facebook AI Research Sorbonne Universités glample@fb.com Myle Ott Facebook AI Research myleott@fb.com

Alexis Conneau Facebook AI Research Université Le Mans aconneau@fb.com

Ludovic Denoyer<sup>†</sup> Sorbonne Universités ludovic.denoyer@lip6.fr

Marc'Aurelio Ranzato Facebook AI Research ranzato@fb.com

#### An Effective Approach to Unsupervised Machine Translation



- First proposed by Lample et al. (2018) proposed similar ideas with SMT, and a joint approach by training NMT on SMT outputs
- Improved by Artetxe et al (2019) with better tuning of the SMT model and a gradual mixing of SMT and NMT backtranslations
- Pure SMT systems perform better than pure NMT systems, yet the best results are obtained by initializing an NMT system with an SMT system

#### **Phrase-Based & Neural Unsupervised Machine Translation**

Guillaume Lample<sup>†</sup> Facebook AI Research Sorbonne Universités glample@fb.com Myle Ott Facebook AI Research myleott@fb.com

Alexis Conneau Facebook AI Research Université Le Mans aconneau@fb.com

Ludovic Denoyer<sup>†</sup> Sorbonne Universités ludovic.denoyer@lip6.fr Marc'Aurelio Ranzato Facebook AI Research ranzato@fb.com

#### An Effective Approach to Unsupervised Machine Translation



• Two recent works (2020) analyzed what are the conditions required to make unsupervised NMT useful

• Two recent works (2020) analyzed what are the conditions required to make unsupervised NMT useful

#### When Does Unsupervised Machine Translation Work?

Kelly Marchisio and Kevin Duh and Philipp Koehn Johns Hopkins University

kmarc@jhu.edu, kevinduh@cs.jhu.edu, phi@jhu.edu



• Two recent works (2020) analyzed what are the conditions required to make unsupervised NMT useful

#### When Does Unsupervised Machine Translation Work?

Kelly Marchisio and Kevin Duh and Philipp Koehn Johns Hopkins University

kmarc@jhu.edu, kevinduh@cs.jhu.edu, phi@jhu.edu

#### When and Why is Unsupervised Neural Machine Translation Useless?

Miguel Graça<sup>†</sup> Hermann Ney Yunsu Kim

Human Language Technology and Pattern Recognition Group RWTH Aachen University, Aachen, Germany {surname}@cs.rwth-aachen.de



- Two recent works (2020) analyzed what are the conditions required to make unsupervised NMT useful
- Both found that:

#### When Does Unsupervised Machine Translation Work?

Kelly Marchisio and Kevin Duh and Philipp Koehn Johns Hopkins University

kmarc@jhu.edu, kevinduh@cs.jhu.edu, phi@jhu.edu

#### When and Why is Unsupervised Neural Machine Translation Useless?

Miguel Graça<sup>†</sup> Hermann Ney Yunsu Kim

Human Language Technology and Pattern Recognition Group RWTH Aachen University, Aachen, Germany {surname}@cs.rwth-aachen.de



- Two recent works (2020) analyzed what are the conditions required to make unsupervised NMT useful
- Both found that:
  - A critical condition is having the data drawn from a **similar domain**

#### When Does Unsupervised Machine Translation Work?

#### Kelly Marchisio and Kevin Duh and Philipp Koehn Johns Hopkins University

kmarc@jhu.edu, kevinduh@cs.jhu.edu, phi@jhu.edu

#### When and Why is Unsupervised Neural Machine Translation Useless?

Miguel Graça<sup>†</sup> Hermann Ney Yunsu Kim

Human Language Technology and Pattern Recognition Group RWTH Aachen University, Aachen, Germany {surname}@cs.rwth-aachen.de

Domain	Domain	BLEU [%]						
(en)	(de/ru)	de-en	en-de	ru-en	en-ru			
Newswire	Newswire	23.3	19.9	11.9	9.3			
	Politics	11.5	12.2	2.3	2.5			
	Random	18.4	16.4	6.9	6.1			



- Two recent works (2020) analyzed what are the conditions required to make unsupervised NMT useful
- Both found that:
  - A critical condition is having the data drawn from a **similar domain**
  - The **pretraining quality** has a strong effect on the final model





Domain	Domain	BLEU [%]						
(en)	(de/ru)	de-en	en-de	ru-en	en-ru			
	Newswire	23.3	19.9	11.9	9.3			
Newswire	Politics	11.5	12.2	2.3	2.5			
	Random	18.4	16.4	6.9	6.1			

### • Unsupervised NMT is possible!



- Unsupervised NMT is possible!
  - Cross-lingual embeddings



- Unsupervised NMT is possible!
  - Cross-lingual embeddings
  - Iterative Back-Translation



- Unsupervised NMT is possible!
  - Cross-lingual embeddings
  - Iterative Back-Translation
- SOTA "Hybrid": SMT bootstrapping, then NMT

Target Source Corpus Corpus Trg Emb Src Emb Space Space Cross-Lingual Emb Space Phrase Phrase Table Table G Tuning with Tuning with SMT System SMT System unsup. MERT unsup. MERT Iterative Iterative backtranslation G > backtranslation SMT System SMT System and tuning with and tuning with MERT MERT -----NMT System NMT System



- Unsupervised NMT is possible!
  - Cross-lingual embeddings
  - Iterative Back-Translation
- SOTA "Hybrid": SMT bootstrapping, then NMT
- When does it work domains, pretraining matters!







