

Objective

Neural Machine Translation (NMT)

Blunsom et al. (2013)

- Motivation:** Handle data from multiple domains by sharing knowledge (Haddow et al. 2012).
- Challenges:** Enforcing knowledge sharing lacks adaptivity to each individual domain.
- Example:** Failure to handle word-level ambiguity across domains: The word "articles" has different meanings in laws and media domains.

Laws	"Article 37 The freedom of marriage ..." "第三十七条:婚姻的自由..."
Media	"... working on an article about the poems ..." "... 为了一篇诗的文章 ..."

Question: How to adaptively capture domain-shared and domain-specific knowledge to improve multi-domain NMT?

Background

Recurrent Network based Encoder-Decoder:

- Computationally Expensive (Recursive Nature)
- Fail to Capture Long-term Dependency
- Various Training Issues (e.g. Gradient Exploding/Vanishing)

Transformer Models (Vaswani et al. 2017)

- Feedforward Network based Encoder-Decoder:

Full Attention, Masked Attention

- Attention (Bahdanau et al. 2014):

Linear Combination, Output

- Multi-Head Attention:

Q, K, V, Scaled Dot-Product Attention, Concat, Point-Wise Linear

Benefits of Transformer Models:

- High Efficiency (Parallel and Feedforward Structures)
- Capture Long-term Dependency (Attention Module)
- Enable Deeper Representation Learning

Proposed Method

Word-Level Domain Proportion

$$D(x) = (1 - \epsilon) \cdot \text{softmax}(Rx) + \epsilon/k$$

- $\epsilon \in (0, 1)$: Smoothing parameter.
- k : Number of domains
- x : Word vector.
- R : Weight matrix of the softmax layer.

Word Level Adaptive Domain Mixing

For a given word vector x_{in} , we take the weighted averaging of the point-wise linear transformations from different domains based on the domain proportion.

$$x_{out} = Wx_{in} \Rightarrow x_{out} = \sum_{j=1}^k D_j(x_{in})W_jx_{in}$$

- W_j : Weight matrix for the j -th domain.
- $D_j(x_{in})$: The proportion of x_{in} for the j -th domain.

Point-wise Linear, Word Level Mixing

Layer-wise Domain Mixing

We apply the word-level domain mixing to the transformer at all layers. Note that the domain proportion is **different** at each attention layer.

Output probabilities, Softmax, Linear, Domain Mixing, Feed Forward, Multihead Attention, Encoder Output, Output Embedding, Outputs (shifted right), Input Embedding, Inputs

Training

$$\min_{\mathcal{H}, \mathcal{D}, \mathcal{F}} L^* = L_{gen}(\mathcal{H}, \mathcal{D}, \mathcal{F}) + L_{mix}(\mathcal{H}, \mathcal{D}, \mathcal{F})$$

- \mathcal{H} : Encoder Module
- \mathcal{F} : Decoder Module
- \mathcal{D} : Domain Proportion
- L_{gen} : Cross Entropy Loss for Translation
- L_{mix} : Cross Entropy Loss over the Smoothed Domain (Hard) Labels of Words
- Training Algorithm: ADAM + Layer Normalization + Step Size Warmup/Annealing

L_{mix} , L_{gen} , Forward, Backward, Input

Experiment – Domain Proportion

Domain Proportion Visualization:

TED Domain (white) vs. Medical Domain (black)

Decoder, Shifted Output, Encoder, Input

Word-Level Analysis:

Top Layers: The phrase is well understood and do not need to borrow knowledge. Ending has little semantic meaning and is shared across domains.

Keep borrowing knowledge from medical domain. Sharing knowledge in Layer 2: Phrase, Layer 1: Single word.

Borrowing little knowledge.

Simple Word, Complicated Word, More common in TED.

Histograms of the Domain Proportions

Within each histogram, 0 means Medical domain, and 1 means TED domain.

Layer-1, 2, 3, 4, 5, 6, Encoder, Decoder

Experiment – Translation

Training Perplexity:

Perplexity, Epoches, News+TED, MTL, AdvL, PAdvL, WDC w/ WL, Mixing: Encoder, Mixing: E/DC

Testing BLEU Scores:

- English to German
- English to French
- Chinese to English

Method	News		TED		TED		Medical	
	News	TED	TED	Medical	TED	Medical	TED	Medical
Direct Training								
News	26.09	6.15			28.22	7.32		
TED	4.90	29.09			7.03	53.73		
News + TED	26.06	28.11			39.21	53.40		
Embedding based Methods								
MTL	26.90	29.27			39.14	53.37		
AdvL	25.68	27.46			39.54	53.46		
PAdvL	27.06	29.49			39.56	53.23		
WDC + WL	27.25	29.43			39.79	53.85		
Our Domain Mixing Methods								
Encoder	27.78	30.30			40.30	54.05		
Encoder + WL	27.67	30.11			40.43	54.14		
E/DC	27.58	30.33			40.52	54.28		
E/DC + WL	27.55	30.22			40.60	54.39		

Method	Laws		News		Speech		Thesis	
	Laws	News	Speech	Thesis	Laws	News	Speech	Thesis
Direct Training								
Laws	51.98	3.80	2.38	2.64				
News	6.88	31.99	8.12	4.17				
Speech	3.33	4.90	18.63	3.08				
Thesis	5.90	5.55	4.77	11.06				
Mixed	48.87	26.92	16.38	12.09				
Embedding based Methods								
MTL	49.14	27.15	16.34	11.80				
AdvL	48.93	26.51	16.18	12.08				
PAdvL	48.72	27.07	15.93	12.23				
WDC + WL	42.16	25.81	15.29	10.14				
Our Domain Mixing Methods								
Encoder	50.21	27.94	16.85	12.03				
Encoder + WL	50.11	27.48	16.79	11.93				
E/DC	50.64	28.48	17.41	11.71				
E/DC + WL	50.04	28.17	17.60	11.59				