



Multi-Round Transfer Learning for Low-Resource NMT Using Multiple High-Resource Languages

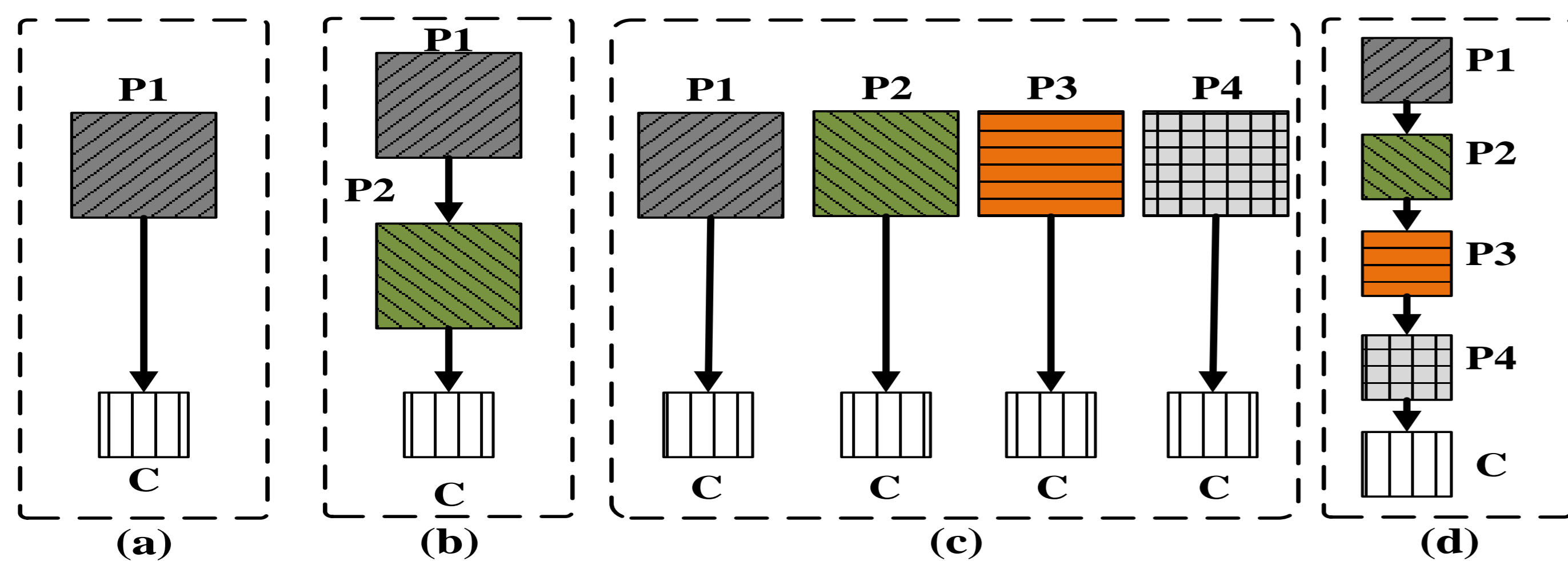


Mieradilijiang Maimaiti¹, Yang Liu¹, Huanbo Luan¹, Maosong Sun¹

¹Department of Computer Science and Technology, Tsinghua University

Motivation and Background

- NMT has made remarkable progress in recent years, but the performance of NMT suffers from a data sparsity problem since large-scale parallel corpora are only readily available for high-resource languages (HRLs).
- Transfer learning (TL) has been used widely in low-resource languages machine translation; while TL is becoming one of the vital directions in low-resource (LR) NMT.



- However, leveraging the original TL to LR models is neither able to make full use of highly related multiple HRLs nor receive different parameters from the same parents.
- To address this issue, we present a language-independent multi-round transfer learning (MRTL) which aims to exploit HRLs effectively. Besides, to reduce the differences between HRLs and LRLs at the character level, we introduce a unified transliteration method for various language families.

Methodology

Main Idea

- we aim to deal with the problem of how to make full use of these corpora of highly related **multiple languages**, to increase the translation quality of the child model.
- Increase the similar even identical words between parent and child language by using **unified transliteration method**.

Multi round fine-tuning

- The original TL transfers parameters of parent model into child model.

$$\theta_{L_3} = \{e_{L_3}, W, e_{L_3}\} \quad \hat{\theta}_{L_3 \rightarrow L_2} = \operatorname{argmax}_{\theta_{L_3 \rightarrow L_2}} \{L(D_{L_3 \rightarrow L_2}, \theta_{L_3 \rightarrow L_2})\} \quad \theta_{L_1 \rightarrow L_2} = f(\hat{\theta}_{L_3 \rightarrow L_2})$$

$$\theta_{L_3 \rightarrow L_2} = f(\hat{\theta}_{L_3 \rightarrow L_2}) \quad \hat{\theta}_{L_3 \rightarrow L_2} = \operatorname{argmax}_{\theta_{L_3 \rightarrow L_2}} \{L(D_{L_3 \rightarrow L_2}, \theta_{L_3 \rightarrow L_2})\} \quad \theta_{L_1 \rightarrow L_2} = f(\hat{\theta}_{L_3 \rightarrow L_2})$$

- The **central idea** of our proposed MRTL is to encourage the child model receive more information from different parent models.

$$\theta_{L_4 \rightarrow L_2} = f(\hat{\theta}_{L_3 \rightarrow L_2}) \quad \hat{\theta}_{L_4 \rightarrow L_2} = \operatorname{argmax}_{\theta_{L_4 \rightarrow L_2}} \{L(D_{L_4 \rightarrow L_2}, \theta_{L_4 \rightarrow L_2})\} \quad \theta_{L_1 \rightarrow L_2} = f(\hat{\theta}_{L_{k+1} \rightarrow L_2})$$

$$\theta_{L_{k+1} \rightarrow L_2} = f(\hat{\theta}_{L_{k+1} \rightarrow L_2}) \quad \hat{\theta}_{L_{k+1} \rightarrow L_2} = \operatorname{argmax}_{\theta_{L_{k+1} \rightarrow L_2}} \{L(D_{L_{k+1} \rightarrow L_2}, \theta_{L_{k+1} \rightarrow L_2})\}$$

Unified transliteration

Algorithm 1: Unified Transliteration Method

```

Input: the source side monolingual sentences  $D_{sm} = \{x_{sm}^m\}_{m=1}^M$  of parent (child).
Output: the transliterated word sequence in current sentences  $D'_{sm}$ .
/* Initialize the variables. */
 $Current_t \leftarrow$  the word sequence in current line among  $D_{sm}$ ;
 $Output_c \leftarrow$  the transliterated word sequence in current line among  $D'_{sm}$ ;
Read the source side monolingual sentences  $D_{sm}$  of parent (child);
for each line in  $D_{sm}$  do
  /* the current line should be decoded as 'utf-8'. */
   $Current_t \leftarrow$  each line.decode('utf-8');
  /* split the current line with white space and save them as a list. */
   $Current_t \leftarrow Current_t.strip().split()$ ;
  for each word in  $Current_t$  do
    for each char in each word do
      /* check each char from the manually prepared mapping table. */
      each  $char_{latin} \leftarrow$  each char;
    end
    /* check the each word if contains same repeated char continually. */
    if Is Contain repeated char in each  $char_{latin}$  then
      compare the length of  $Current_t$  and the length of Latinized  $Current_t$ ;
      remove repeated each  $char_{latin}$  from each word;
    end
    convert them into unified form sequentially;
    each word  $\leftarrow$  sums of unified chars after removing repeated chars;
    final word  $\leftarrow$  each word; /* reserve the final word. */
  end
 $Output_c \leftarrow$  the joint sequence of final words with white-space
end

```

Experiments

Datasets

Language features							Characteristics of corpora							
Language	Family	Group	Branch	Order	Unit	Inflection	Languages	Train	Dev	Test	Source	Target		
											Vocab.	# Word	Vocab.	# Word
Arabic (Ar)	Hamito-Semitic	Semitic	South	VSO	Word	High	Ar → Ch	5.1M	2.0K	2.0K	1.0M	32.2M	0.5M	37.4M
Farsi (Fa)	Indo-European	Indic	West	SOV	Word	Moderate	Fa → Ch	1.4M	2.0K	1.0K	0.2M	10.4M	0.2M	10.0M
Urdu (Ur)	Indo-European	Iranian	Iranian	SOV	Word	Moderate	Ur → Ch	78.0K	1.0K	1.0K	17.6K	2.6M	12.7K	2.4M
Finnish (Fi)	Uralic	Finno-Ugric	Finnish	SVO	Word	Moderate	Fi → Ch	2.8M	2.0K	1.0K	0.7M	18.4M	0.3M	23.1M
Hungarian (Hu)	Uralic	Finno-Ugric	Ugric	SVO	Word	Moderate	Hu → Ch	4.1M	2.0K	1.0K	1.0M	30.4M	0.5M	32.5M
Turkish (Tr)	Altaic	Turkic	Oghuz	SOV	Word	Moderate	Tr → Ch	4.4M	2.0K	1.0K	0.7M	30.6M	0.5M	35.9M
Uyghur (Uy)	Altaic	Turkic	Qarluq	SOV	Word	Moderate	Uy → Ch	46.3K	1.0K	1.0K	73.5K	1.1M	42.1K	11.2M
Chinese (Ch)	Sino-Tibetan	Chinese	Sinitic	SVO	Character	Light								

The "Vocab" and "# Word" represent vocabulary (word type) and word token, respectively.

Unified transliterations

Method	Round	Parent	Child	BLEU
TRANSFORMER	R=0	N/A		28.28
		Ur → Ch		10.29
MRTL (Original)	R=1	Fa → Ch	Uy → Ch	28.83
		Ar → Ch		30.64 ⁺⁺
MRTL (Unified)	R=1	Ur → Ch		10.93 [*]
		Fa → Ch		29.96 ⁺⁺⁺
		Ar → Ch		31.64 ⁺⁺⁺

Effect of parent models

Method	Parent	Child	BLEU
TRANSFORMER	N/A		28.28
	Ur → Ch		10.93
MRTL (R=1)	Fa → Ch		29.96 ⁺⁺
	Fi → Ch	Uy → Ch	30.85 ⁺⁺
	Tr → Ch (2.4M)		30.88 ⁺⁺⁺
	Ar → Ch		31.64 ⁺⁺⁺
	Hu → Ch		32.41 ⁺⁺⁺
	Tr → Ch (4.4M)		32.74 ⁺⁺⁺

Parent language selection

Different language family

MRTL	Parent	Family	Domain	Size	Child	BLEU
R=0	N/A	Altaic	CLDC	46.3K		28.28
R=1	Hu → Ch	Uralic	Open Subtitle	4.1M	Uy → Ch	32.41 ⁺⁺
	Tr → Ch	Altaic				32.58 ⁺⁺

Different domain

MRTL	Parent	Family	Domain	Size	Child	BLEU
R=0	N/A	Altaic	CLDC	46.3K		28.28
R=1	Ur → Ch	Indo-European	Tanzil	78.0K	Uy → Ch	10.93
	Fa → Ch		Open Subtitle			24.27

Different corpus size

MRTL	Parent	Family	Domain	Size	Child	BLEU
R=0	N/A	Altaic	CLDC	46.3K		28.28
R=1	Fi → Ch	Uralic	Open Subtitle	2.8M	Uy → Ch	30.85 ⁺⁺
	Hu → Ch			4.1M		32.41 ⁺⁺

Effect of MRTL method

Method	Round	Parent	Child	BLEU
TRANSFORMER	R=0	N/A		28.28
MANY-TO-ONE	R=1	Tr (4.4M) → Ch		32.43 ⁺⁺
		Tr (4.4M), (2.4M) → Ch		32.03 ⁺⁺
MRTL	R=2	Tr (4.4M), (2.4M) → Ch	Uy → Ch	32.54 ⁺⁺
	R=3	Tr (4.4M), (2.4M), Fi → Ch		33.54 ⁺⁺⁺
		Tr (4.4M), (2.4M), Fi, Hu → Ch		33.66 ⁺⁺⁺
	R=4	Ar (Unified), Tr (4.4M), Hu, Fi → Ch		33.73 ⁺⁺⁺
		Tr (4.4M), Ar (Unified), Hu, Fi → Ch		33.91 ⁺⁺⁺

Examples

Method	Translation result
Source	<i>muvaپیق sürük İçide yolğa qoyu@lmisa dölet heqsiz yolğa qoysa bolidu .</i>
Reference	<i>zai heli qixian nei meiyou shishi de , guo jia keyi wuchang shishi .</i> 在合理期限内没有实施的，国家可以无偿实施。
TRANSFORMER	<i>shidang xianqi shishi .</i> 适当限期实施。
MANY-TO-ONE	<i>zai shidang qixian nei bu neng shixing guo jia mianfei shishi .</i> 在适当期限内不能实行国家免费实施。
R = 1	<i>zai shidang qixian nei bu neng you mianfei shishi guo jia .</i> 在适当期限内不能有免费实施国家。
R = 2	<i>zai heli qixian nei shishi xi ze , guo jia keyi textbfmianfei shishi .</i> 在合理期间内实施细则，国家可以免费实施。
R = 3	<i>zai heli qixian nei wei shixing guo jia ke mianfei shishi .</i> 在合理期间内未实行国家可免费实施。
R = 4	<i>dui heli qixian nei wei shixing de , guo jia keyi wuchang shishi .</i> 对合理期间内未实行的，国家可以无偿实施。

Conclusion

- We address the drawbacks of TL, which exploits only one parent to optimize the child model at a time.
- We mitigate the gap between parent and child language pairs at the character level.
- We achieve transparency in network architectures, as well as in our method for neural network architecture.
- We observe meaningful discovery by sharing both source side and target side embeddings of parent models.