

NANYANG

UNIVERSITY

SINGAPORE

TECHNOLOGICAL



Overview

- We present an approach to this *unpaired image captioning* problem by language pivoting.
- Our method can effectively capture the characteristics of an image captioner from the *pivot language (Chinese)* and align it to the target language (English) using another pivot-target (Chinese-English) parallel corpus.
- Quantitative comparisons against several baseline approaches demonstrate the effectiveness of our method.

Challenges

- The task of image caption generation has so far been explored only in English since most available datasets are in English. However, only 5.5% of the world population are native speakers of English (source from Wikipedia).
- ²In many applications and languages, the large-scale annotations are not readily available, and are expensive and slow to acquire.
- **Getting** $image \rightarrow caption paired data for each target language is$ expensive, but, language parallel data is relatively easier to get.



Motivation

- Our pivot-based captioning framework learns to:
- Generate captions for $image \rightarrow pivot$ language in domain A.
- Translate sentences for $pivot \rightarrow target$ language in domain B.
- Match it to the target language with caption style in *domain* C.
- The two proposed connection terms $(\mathcal{R}_{i \to y} \text{ and } \mathcal{R}_{x \to \hat{y}})$ bridge the gap between the pivot language in *domain* A and *domain* B and the gap between the target language in *domain* B and *domain* C, respectively.

Unpaired Image Captioning by Language Pivoting

Jiuxiang Gu¹, Shafiq Joty¹, Jianfei Cai¹, Gang Wang² 1 Nanyang Technological University, 2 Alibaba Al Labs

Pivot-based image captioning framework



• We use the pivot language x to learn the mapping: $i \xrightarrow{\theta_{i \to x}} x \xrightarrow{\theta_{x \to y}} y$. Our pivot-based image captioning approach has an image captioning model $P(x|i;\theta_{i\to x})$ to generate a caption x in the pivot language from an image i and a machine translation model $P(y|x; \theta_{x \to y})$ to translate this caption into the target language. Build and a subset of the sentences of the sentences of the sentences of the sentences of the sentences. • We train these components jointly so that they interact with each other.

Learning

1 We first generate the image description x from the global image feature v. We maximize the probability of the ground truth caption words given the image:

$$\tilde{\theta}_{i \to x} = \arg \max_{\theta_{i \to x}} \left\{ \mathcal{L}_{i \to x} \right\} = \arg \max_{\theta_{i \to x}} \left\{ \sum_{\substack{\Sigma \\ n_i = 0}}^{N_i - 1} \frac{M^{(n_i)} - 1}{t = 0} \log P_x(x_t^{(n_i)} | x_{0:t-1}^{(n_i)}, i^{(n_i)}; \theta_{i \to x}) \right\}$$
(1)
second model is to learn to translate the Chinese sentence x to the English

²The second model is to learn to translate the Chinese sentence x to the English sentence y, where the maximum-likelihood training objective of the model can be expressed as:

$$\tilde{\theta}_{x \to y} = \arg \max_{\theta_{x \to y}} \left\{ \mathcal{L}_{x \to y} \right\} = \arg \max_{\theta_{x \to y}} \left\{ \sum_{\substack{\Sigma \\ n_x = 0 \ t = 0}}^{N_x - 1} \log P_y(y_t^{(n_x)} | y_{0:t-1}^{(n_x)}; x^{(n_x)}; \theta_{x \to y}) \right\}$$
(2)
3 Connecting *Image-to-Pivot* and *Pivot-to-Target*.

$$\mathcal{R}_{i \to u}(\theta_{i \to w}^{w_x}, \theta_{x \to u}^{w_x}) = - \qquad \Sigma \qquad ||\theta_{i \to w}^{w_x} - \theta_{x \to u}^{w_x}||_2$$

$$\mathcal{R}_{i \to y}(\theta_{i \to x}^{w_x}, \theta_{x \to y}^{w_x}) = -\sum_{\substack{w_x \in \mathcal{V}_{i \to x}^x \cap \mathcal{V}_{x \to y}^x \\ w_y \in \mathcal{V}_{i \to x}^y \cap \mathcal{V}_{y \to y}^y}} ||\theta_{i \to x}^{w_x} - \theta_{x \to y}^{w_x}||_2$$
(3)
$$\mathcal{R}_{x \to \hat{y}}(\theta_{x \to y}^{w_y}, \theta_{\hat{y} \to \hat{y}}^{w_y}) = -\sum_{\substack{w_y \in \mathcal{V}_{x \to y}^y \cap \mathcal{V}_{\hat{y} \to \hat{y}}^y \\ w_y \in \mathcal{V}_{x \to y}^y \cap \mathcal{V}_{\hat{y} \to \hat{y}}^y}} ||\theta_{x \to y}^{w_y} - \theta_{\hat{y} \to \hat{y}}^{w_y}||_2$$
(4)

• where $\theta_{i \to x}^{w_x}$ is already a learned model and kept fixed during adaptation.

• Our goal is to find a set of source-to-target model parameters that maximizes the training objective:

$$\mathcal{J}_{i \to x, x \to y, y \to \hat{y}} = \mathcal{L}_{i \to x} + \mathcal{L}_{x \to y} + \mathcal{L}_{\hat{y} \to \hat{y}} + \lambda \mathcal{R}_{i \to x, x \to y, y \to \hat{y}}$$
(5)

$$\mathcal{R}_{i \to x, x \to y, y \to \hat{y}} = \mathcal{R}_{i \to y}(\theta_{i \to x}^{w_x}, \theta_{x \to y}^{w_x}) + \mathcal{R}_{x \to \hat{y}}(\theta_{x \to y}^{w_y}, \theta_{\hat{y} \to \hat{y}}^{w_y})$$
(6)

6 During inference, given an unseen image i to be described, we use the joint decoder:

$$y \sim \arg\max_{y} \{P(y|i; \theta_{i \to x}, \theta_{x \to y})\}$$
(7)

Implementation



(a) Image captioning model (b) Machine translation model

(c) Autoencoder

We choose the two independent datasets used from AI Challenger (AIC): AIC Image Chinese Captioning (AIC-ICC) and AIC Chinese-English Machine Translation (AIC-MT), as the training datasets, while using MSCOCO and Flickr30K as the test datasets.

	Dataset	Lang.	Sourc	ce	Target		
			# Image/Sent.	Vocab. Size	# Sent.	Vocab. Size	
Training	AIC-ICC	$\operatorname{im} \to \operatorname{zh}$	240K	_	1,200K	4,461	
	AIC-MT	$zh \rightarrow en$	10,000 K	50,004	10,000 K	50,004	
Testing	MSCOCO	$\operatorname{im} \to \operatorname{en}$	123K		615K	9,487	
	Flickr30K	$ \mathrm{im} \rightarrow \mathrm{en} $	30K		150K	7,000	



ECCV2018European Conference on Computer Vision 8 – 14 September 2018 I Munich, Germany

Word clouds of datasets



With





(a) AIC-MT (English

(c) AIC-MT (Chinese)

(d) AIC-ICC (Chinese)

Results



Results of unpaired image-to-English captioning on MSCOCO 5K test split.

Approach	Lang.	B@1	B@2	B@3	B@4	Μ	CIDEr
	MSCC	CO					
$i2t_{im \rightarrow en}$ (Upper bound, XE Loss)	en	73.2	56.3	42.0	31.2	25.3	95.1
FC-2K [42] (ResNet101, XE Loss)	en		_		29.6	25.2	94.0
$\overline{\mathrm{i}2\mathrm{t}_{\mathrm{im}\to\mathrm{zh}} + \mathrm{nm}\mathrm{t}_{\mathrm{zh}\to\mathrm{en}} \left(\mathcal{R}_{i\to x,x\to y,y\to\hat{y}}\right)}$	en	46.2	24.0	11.2	5.4	13.2	17.7
$i2t_{im \to zh} + nmt_{zh \to en} (\mathcal{R}_{i \to x, x \to y})$	en	45.5	23.6	11.0	5.3	13.1	17.3
$i2t_{im \rightarrow zh} + nmt_{zh \rightarrow en}$ (Lower bound)	en	42.0	20.6	9.5	3.9	12.0	12.3
$i2t_{im \rightarrow zh}$ + Online Google Translation	en	42.2	21.8	10.7	5.3	14.5	17.0
	Flickr	30K					
$i2t_{im \rightarrow en}$ (Upper bound, XE Loss)	en	63.1	43.8	30.2	20.7	17.7	40.1
$\overline{\mathrm{i2t_{im}}_{zh} + \mathrm{nmt_{zh}}_{en} (\mathcal{R}_{i \to x, x \to y, y \to \hat{y}})}$	en	49.7	27.8	14.8	7.9	13.6	16.2
$i2t_{im \to zh} + nmt_{zh \to en} (\mathcal{R}_{i \to x, x \to y})$	en	48.7	26.1	12.8	6.4	13.0	14.9
$i2t_{im \rightarrow zh} + nmt_{zh \rightarrow en}$ (Lower bound)	en	45.9	25.2	13.1	6.9	12.5	13.9
$i2t_{im \rightarrow zh}$ + Online Google Translation	en	46.2	25.4	13.9	7.7	14.4	15.8

Evaluation results of user assessment on MSCOCO 1.2K test split.

Approach	$i2t_{im \to zh} + nmt_{zh \to en} (\mathcal{R}_{i \to x, x \to y, y \to \hat{y}})$	Upper Bound	Ground-Truth
Relevant	3.81	3.99	4.68
Resemble	3.78	4.05	4.48