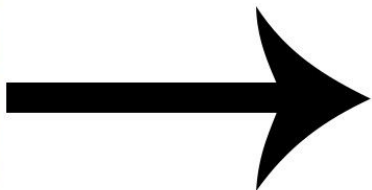


Distinctive Image Captioning: Leveraging Ground Truth Captions in CLIP Guided Reinforcement Learning

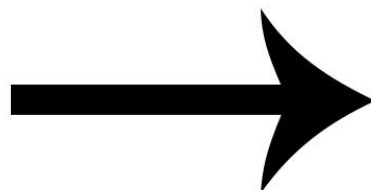
Antoine Chaffin,
Vincent Claveau,
Ewa Kijak



- Language model conditioned on an image
- Create a **powerful cross-modal alignment**^[1]



Language model



A cat on a branch

- Datasets captions only describe most salient objects, common to many images
- Higher word-matching metrics with words common across different images, not specific ones

A couple of dogs standing on a porch



- Datasets captions only describe most salient objects, common to many images
- Higher word-matching metrics with words common across different images, not specific ones

A couple of dogs standing on a porch



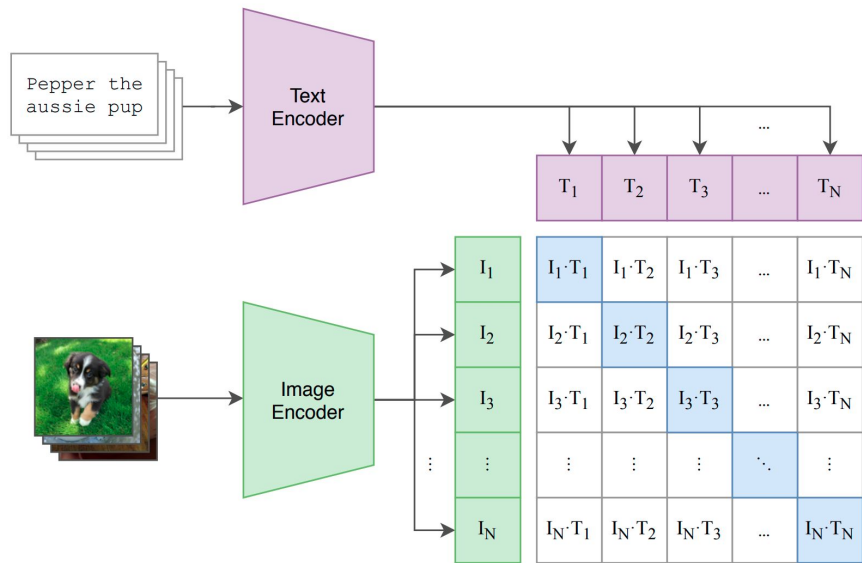
- Fine-grained alignment to describe **this image and only this one**

- Reinforcement learning to optimize cross-modal similarity of the generated caption and the target image
 - A description that can let the retriever identify the image

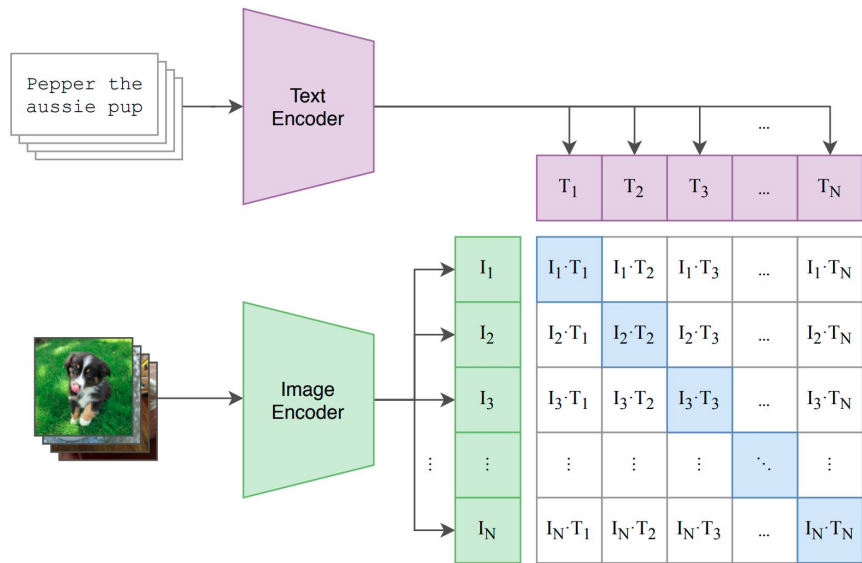


a couple of dogs wearing a santa hat on a porch

- Dual encoder, each projecting a modality separately
 - Similarity using dot product of both representations



- Dual encoder, each projecting a modality separately
 - Similarity using dot product of both representations
- Couple closer than any element in the batch



$$\mathcal{L}_{\text{CLIP}} = \underbrace{\log \frac{e^{\frac{t_c \cdot i_c}{\tau}}}{\sum_{t \in \mathcal{T}} e^{\frac{t \cdot i_c}{\tau}}}}_{\text{image-to-text}} + \underbrace{\log \frac{e^{\frac{t_c \cdot i_c}{\tau}}}{\sum_{i \in \mathcal{I}} e^{\frac{t_c \cdot i}{\tau}}}}_{\text{text-to-image}}$$

- Prevent the model from learning ill-formed solutions



*a close up of two **brown** and **black** dogs wearing a **santa hat** on a **black** and **brown dog** with a **red hat** on a backyard with a fence in the background*

- Prevent the model from learning ill-formed solutions
- Regularization term in the reward
 - KL divergence, CIDEr value, **grammar network**...



*a close up of two **brown** and **black** dogs wearing a **santa hat** on a **black** and **brown dog** with a **red hat** on a backyard with a fence in the background*

$$\nabla_{\theta} L_{\theta}(x) = - \left[\underbrace{\alpha r_{sim}(x)}_{\text{Similarity reward}} + \underbrace{(1 - \alpha) r_{regu}(x)}_{\text{Regularization reward}} \right] \nabla_{\theta} \underbrace{\log p_{\theta}(x)}_{\text{Likelihood}}$$

Sample from the generator
Likelihood

- 3 different contributions to improve CLIP-based RL image captioning
 1. **Discriminator regularization**
 2. **RL objective on ground truth samples**
 3. **Bidirectional contrastive reward**

- 3 different contributions to improve CLIP-based RL image captioning
 1. **Discriminator regularization**
 2. **RL objective on ground truth samples**
 3. **Bidirectional contrastive reward**
- MS COCO dataset
- Trade-off:
 - **Discriminativeness**: recall@k using generated caption (fixed CLIP model)
 - **Writing quality**: BLEU, ROUGE, CIDEr, METEOR and SPICE

- Use generated text discriminator scores as regularization
- Simple MLP using CLIP representations as input

$$\nabla_{\theta} L_{\theta}(x) = - \left[\left(\alpha r_{sim}(x) + (1 - \alpha) r_{regu}(x) \right) \nabla_{\theta} \log p_{\theta}(x) \right]$$

Diagram illustrating the gradient of the loss function $\nabla_{\theta} L_{\theta}(x)$ with annotations:

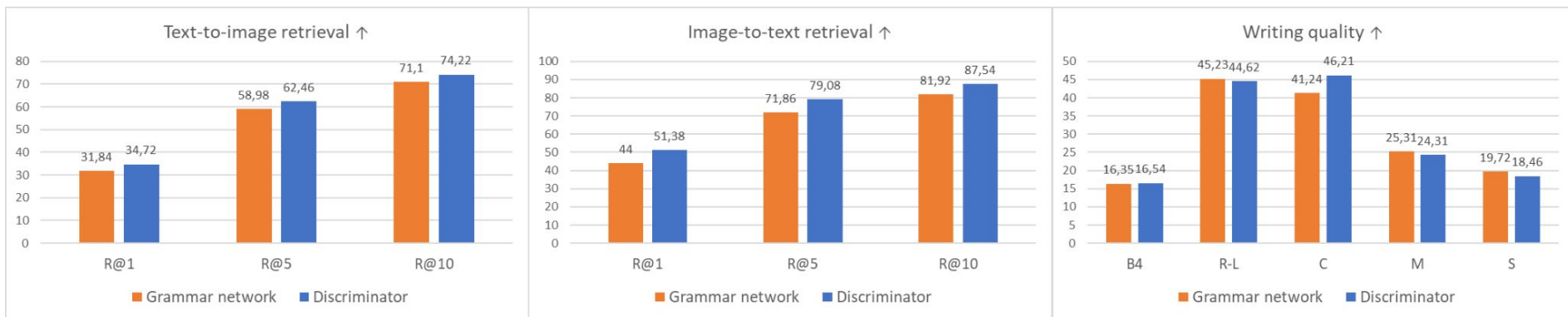
- Sample from the generator**: Points to x in $\nabla_{\theta} L_{\theta}(x)$.
- Similarity reward**: Points to $\alpha r_{sim}(x)$.
- Regularization reward**: Points to $(1 - \alpha) r_{regu}(x)$.
- Likelihood**: Points to $\log p_{\theta}(x)$.

- Use generated text discriminator scores as regularization
- Simple MLP using CLIP representations as input

$$\nabla_{\theta} L_{\theta}(x) = - \left[\left(\alpha r_{sim}(x) + (1 - \alpha) r_{regu}(x) \right) \nabla_{\theta} \log p_{\theta}(x) \right]$$

Similarity reward (points to $\alpha r_{sim}(x)$)
Regularization reward (points to $(1 - \alpha) r_{regu}(x)$)
Sample from the generator (points to x)
Likelihood (points to $\log p_{\theta}(x)$)

- Higher retrieval rate without degrading written quality



- RL learns from high-scoring sequences
- Ground truths are (relatively) good solutions

$$\nabla_{\theta} L_{\theta}(x) = - \left(\alpha r_{sim}(x) + (1 - \alpha) r_{regu}(x) \right) \nabla_{\theta} \log p_{\theta}(x)$$

Diagram illustrating the components of the loss function $L_{\theta}(x)$ and its gradient $\nabla_{\theta} L_{\theta}(x)$:

- Ground truth sample** (x) is the input to the loss function.
- The loss function is composed of two terms: $\alpha r_{sim}(x)$ (Similarity reward) and $(1 - \alpha) r_{regu}(x)$ (Regularization reward).
- The gradient of the loss function is $\nabla_{\theta} L_{\theta}(x)$.
- The gradient is calculated as the negative of the weighted sum of the rewards, multiplied by the gradient of the log-likelihood: $\nabla_{\theta} \log p_{\theta}(x)$ (Likelihood).

- RL learns from high-scoring sequences
- Ground truths are (relatively) good solutions
- Learn to reproduce human-written sequence (TF) but focuses on highly descriptive ones

$$\nabla_{\theta} L_{\theta}(x) = - \left(\alpha r_{sim}(x) + (1 - \alpha) r_{regu}(x) \right) \nabla_{\theta} \log p_{\theta}(x)$$

Similarity reward

Regularization reward

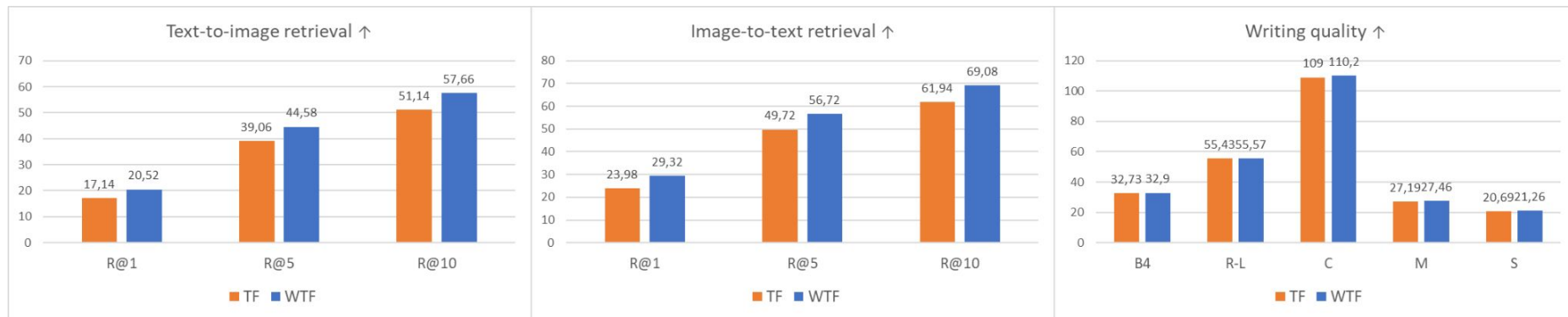
Ground truth sample

Likelihood



- (1) *there is an adult bear that is walking in the forest*
- (2) *picture of an exterior place that looks wonderful.*

- Improve retrieval metrics using only ground truth, without degrading writing quality
- Better regularization objective to couple with traditional RL



- Subtract a baseline to the reward to reduce variance

$$\nabla_{\theta} L_{\theta}(x) = - \left(\overset{\text{Reward}}{r(x)} - \overset{\text{Baseline}}{b} \right) \nabla_{\theta} \log p_{\theta}(x)$$

Sample from the generator

Likelihood

- Subtract a baseline to the reward to reduce variance

$$\nabla_{\theta} L_{\theta}(x) = - \left(\overset{\text{Reward}}{r(x)} - \overset{\text{Baseline}}{b} \right) \nabla_{\theta} \overset{\text{Likelihood}}{\log p_{\theta}(x)}$$

↑ Sample from the generator

1. Use the model itself as a baseline^[1]

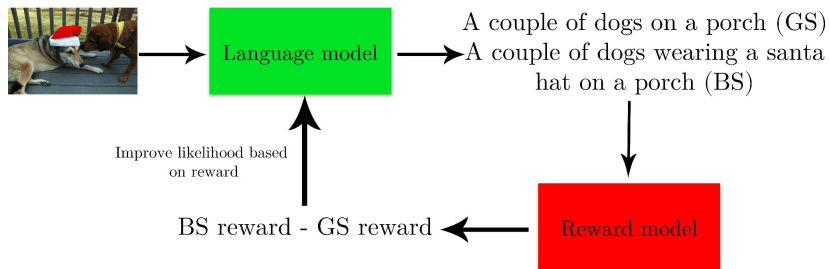


Image-to-text baseline

- Subtract a baseline to the reward to reduce variance

$$\nabla_{\theta} L_{\theta}(x) = - \left(\overset{\text{Reward}}{r(x)} - \overset{\text{Baseline}}{b} \right) \nabla_{\theta} \overset{\text{Likelihood}}{\log p_{\theta}(x)}$$

↑ Sample from the generator ↑ Likelihood

- Use the model itself as a baseline^[1]
- Similarity with other (similar) images^[2]

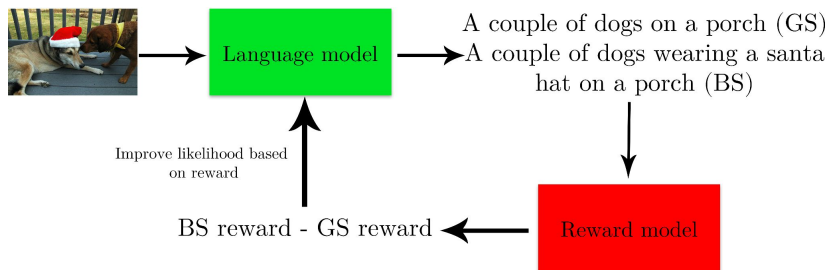
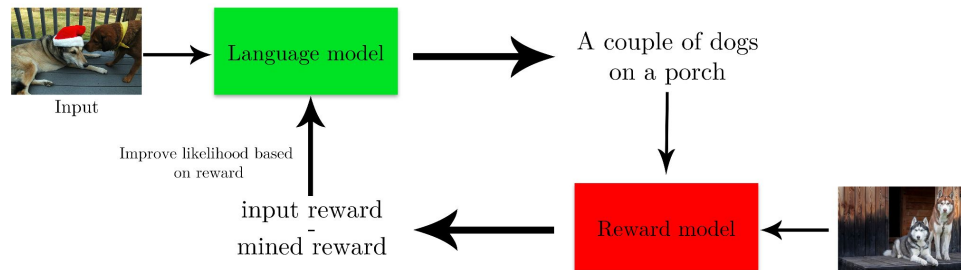


Image-to-text baseline






Text-to-image baseline

[1] Jaemin Cho, Seunghyun Yoon, Ajinkya Kale, Franck Deroncourt, Trung Bui, Mohit Bansal. "Fine-grained Image Captioning with CLIP Reward". 2022




[2] Youyuan Zhang, Jiuniu Wang, Hao Wu, Wenjia Xu. "Distinctive Image Captioning via CLIP Guided Group Optimization". 2022

- Decoupled contrastive loss

		A couple of dogs wearing a santa hat (BS 1)	A couple of dogs on a porch (GS 1)	A couple of dogs wearing a santa hat on a porch (GT 1)	A cat on a branch (BS N)		
Input image 1		$I_1 \cdot T_1$	$I_1 \cdot T_2$	$I_1 \cdot T_3$...	$I_1 \cdot T_N$...
Mined image 1		$I_2 \cdot T_1$	$I_2 \cdot T_2$	$I_2 \cdot T_3$...	$I_2 \cdot T_N$...
		⋮	⋮	⋮	⋮	⋮	⋮
Input image N		$I_N \cdot T_1$	$I_N \cdot T_2$	$I_N \cdot T_3$...	$I_N \cdot T_N$...
		⋮	⋮	⋮	⋮	⋮	⋮

$$r_{bicont}(t_c) = \tau \left(\underbrace{\log \frac{e^{\frac{t_c \cdot i_c}{\tau}}}{\sum_{t \in \mathcal{T} \setminus t_c} e^{\frac{t \cdot i_c}{\tau}}}}_{\text{Image-to-text reward } r_{i2t}(t_c)} + \underbrace{\log \frac{e^{\frac{t_c \cdot i_c}{\tau}}}{\sum_{i \in \mathcal{I} \setminus i_c} e^{\frac{t_c \cdot i}{\tau}}}}_{\text{Text-to-image reward } r_{t2i}(t_c)} \right)$$

- Decoupled contrastive loss
- Closest element in the batch as baseline
- Natively handle both cross-modal directions




	A couple of dogs wearing a santa hat (BS 1)	A couple of dogs (GS 1)	A couple of dogs wearing a santa hat on a porch (GT 1)	...	A cat on a branch (BS N)	...	
Input image 1		$I_1 \cdot T_1$	$I_1 \cdot T_2$	$I_1 \cdot T_3$...	$I_1 \cdot T_N$...
Mined image 1		$I_2 \cdot T_1$	$I_2 \cdot T_2$	$I_2 \cdot T_3$...	$I_2 \cdot T_N$...
	⋮	⋮	⋮	⋮	⋮	⋮	⋮
Input image N		$I_N \cdot T_1$	$I_N \cdot T_2$	$I_N \cdot T_3$...	$I_N \cdot T_N$...
	⋮	⋮	⋮	⋮	⋮	⋮	⋮

$$r_{bicont}(t_c) = \tau \left(\underbrace{\log \frac{e^{\frac{t_c \cdot i_c}{\tau}}}{\sum_{t \in \mathcal{T} \setminus t_c} e^{\frac{t \cdot i_c}{\tau}}}}_{\text{Image-to-text reward } r_{i2t}(t_c)} + \underbrace{\log \frac{e^{\frac{t_c \cdot i_c}{\tau}}}{\sum_{i \in \mathcal{I} \setminus i_c} e^{\frac{t_c \cdot i}{\tau}}}}_{\text{Text-to-image reward } r_{t2i}(t_c)} \right)$$

$$r_{i2t}(t_c) = \tau \left(\log \left(e^{\frac{t_c \cdot i_c}{\tau}} \right) - \log \left(\sum_{t \in \mathcal{T} \setminus t_c} e^{\frac{t \cdot i_c}{\tau}} \right) \right)$$

$$\approx t_c \cdot i_c - \max_{t \in \mathcal{T} \setminus t_c} \{t \cdot i_c\}$$

- Decoupled contrastive loss
- Closest element in the batch as baseline
- Natively handle both cross-modal directions
- The caption is **very descriptive of the image and this image only**

		A couple of dogs wearing a santa hat (BS 1)	A couple of dogs on a porch (GS 1)	A couple of dogs wearing a santa hat on a porch (GT 1)	A cat on a branch (BS N)	
Input image 1		$I_1 \cdot T_1$	$I_1 \cdot T_2$	$I_1 \cdot T_3$...	$I_1 \cdot T_N$
Mined image 1		$I_2 \cdot T_1$	$I_2 \cdot T_2$	$I_2 \cdot T_3$...	$I_2 \cdot T_N$
		⋮	⋮	⋮	⋮	⋮
Input image N		$I_N \cdot T_1$	$I_N \cdot T_2$	$I_N \cdot T_3$...	$I_N \cdot T_N$
		⋮	⋮	⋮	⋮	⋮

$$r_{bicont}(t_c) = \tau \left(\underbrace{\log \frac{e^{\frac{t_c \cdot i_c}{\tau}}}{\sum_{t \in \mathcal{T} \setminus t_c} e^{\frac{t \cdot i_c}{\tau}}}}_{\text{Image-to-text reward } r_{i2t}(t_c)} + \underbrace{\log \frac{e^{\frac{t_c \cdot i_c}{\tau}}}{\sum_{i \in \mathcal{I} \setminus i_c} e^{\frac{t_c \cdot i}{\tau}}}}_{\text{Text-to-image reward } r_{t2i}(t_c)} \right)$$

$$r_{i2t}(t_c) = \tau \left(\log \left(e^{\frac{t_c \cdot i_c}{\tau}} \right) - \log \left(\sum_{t \in \mathcal{T} \setminus t_c} e^{\frac{t \cdot i_c}{\tau}} \right) \right)$$

$$\approx t_c \cdot i_c - \max_{t \in \mathcal{T} \setminus t_c} \{t \cdot i_c\}$$

- Unidirectional image-to-text reward only yield significantly lower text-to-image retrieval results
- Both cross-modal directions are needed for a caption highly descriptive of **this image and this image only**

