

# VLMixer

Unpaired Vision-Language Pre-training via CMC

---

Mose Park

Department of Statistical Data Science  
University of Seoul

Selective. Lab

May 23, 2024

---

# Index

- 1 Problem
  - 2 Approach
  - 3 Details
  - 4 Experiment
-

# Main concepts

```
graph LR; A((Vision-Language Data Augmentation)) --- B((Modality Gap)); B --- C((Loss function));
```

Vision-Language  
Data Augmentation

Modality Gap

Loss function

# 1

## Problem

---

# Problem

Paired



\* image

↔ train traveling down a track in front of road  
\* text

unpaired



\* image

⋯ ? ⋯

- A train traveling down tracks next to lights.
- A blue train is next to a sidewalk on the rails.
- A passenger train pulls into a train station.

\* text

► Obtaining paired image-text data is a **costly** and precise task. → Need **data augmentation**!

# Paired VLP

Learning Obj

single stream  
dual stream

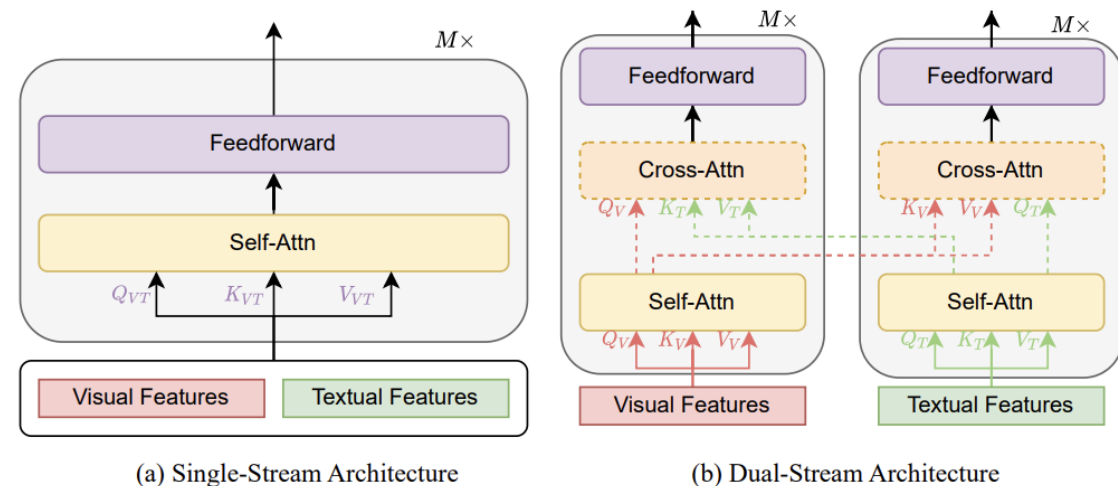
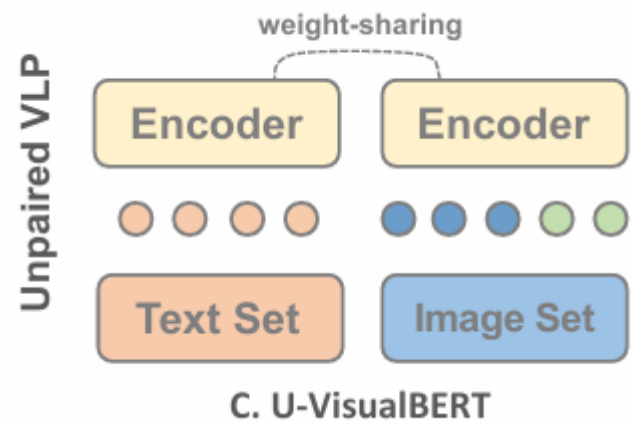


Fig. 1 Illustration of two types of model architectures for VLP.

# Unpaired VLP



\* Same system, Sharing parameter

\* Independent dataset

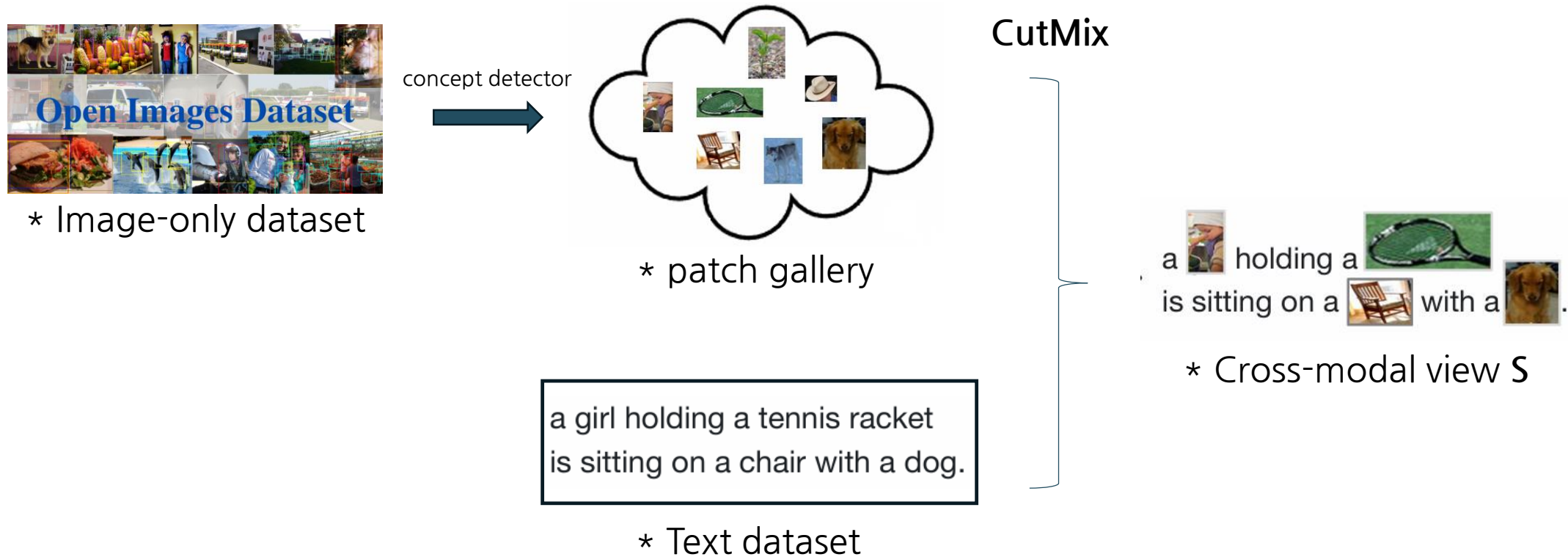
► How does this paper address the modality gap?

# 2

## Approach

---

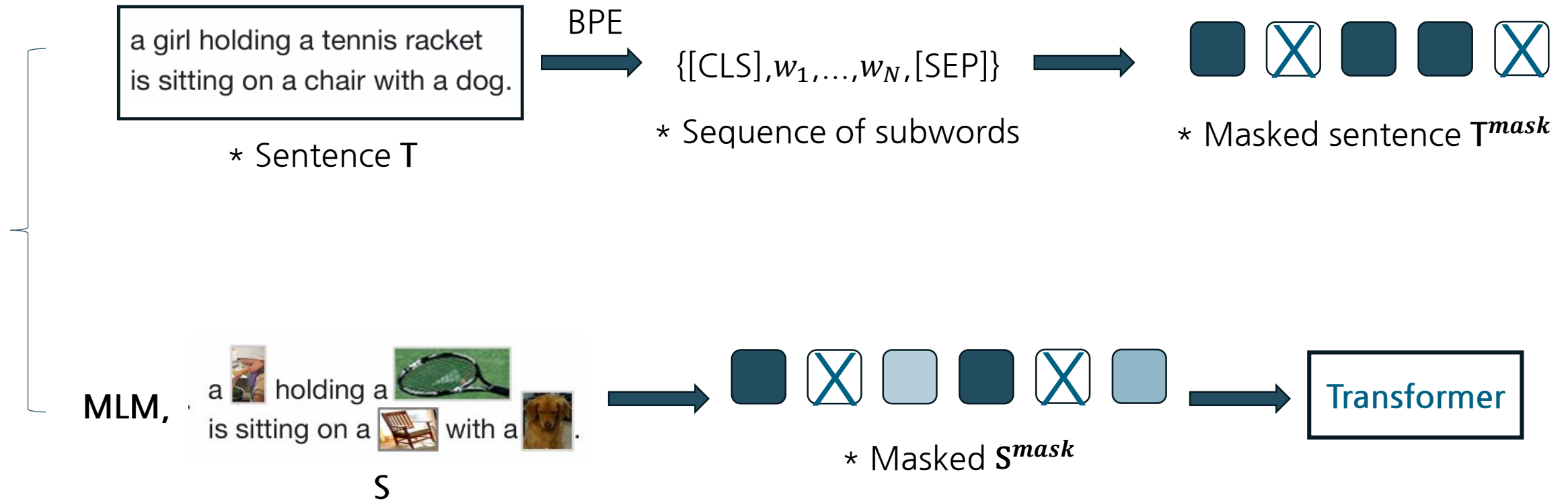
# Approach : Cross-modal CutMix



- ▶ Preserve syntactic and semantic information
- ▶ Introduce visual tokens as the cross-modality noise



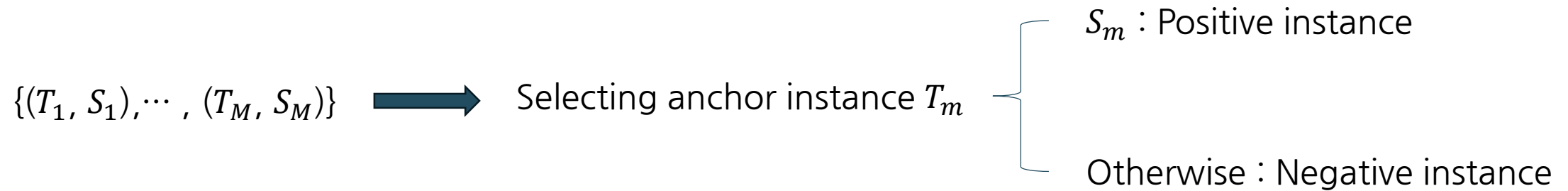
# Approach : VALP



$$\mathcal{L}_{\text{mlm}} = -\mathbb{E}_{\mathbf{T} \sim \mathcal{D}^T} \log(\mathbf{T} | \mathbf{S}^{\text{mask}}) = \text{CE}(\hat{\mathbf{S}}, \mathbf{T})$$

- **Masked language model** exploits multi-modal fusion.
- **Contrastive learning** is to learn cross-modal alignments.

## Contd.

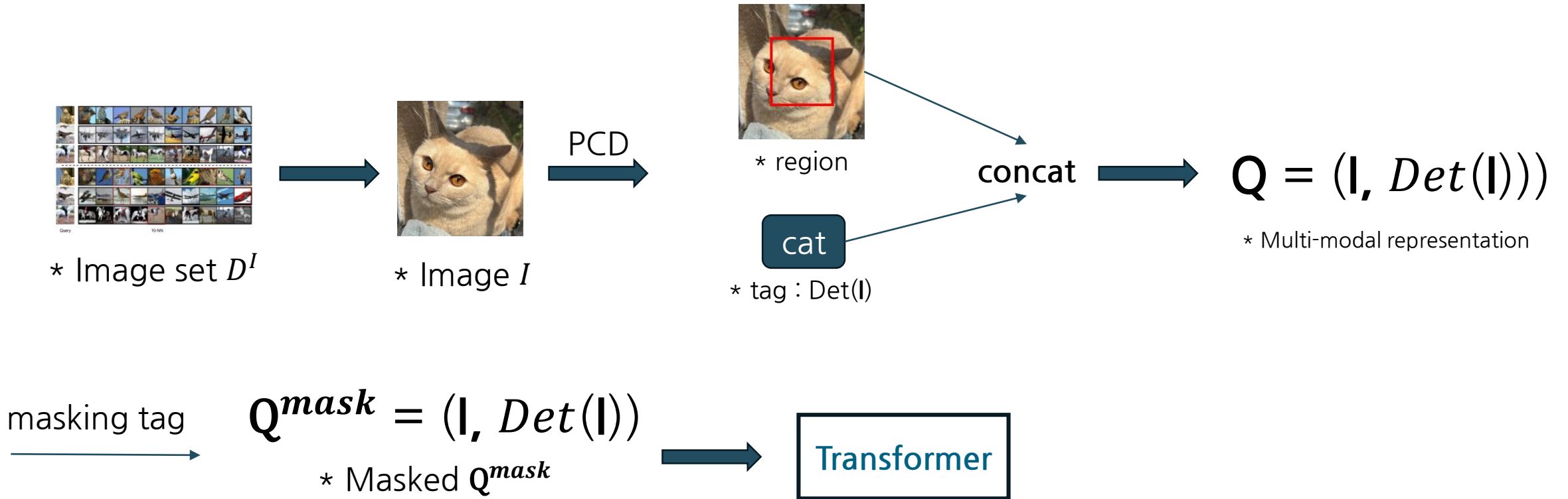


$\longrightarrow$  Contrastive Learning ! (By Eqn (4) Loss)

$$\mathcal{L}_{\text{cl}} = - \sum_{m=1}^M \log \frac{\exp(f(\mathbf{T}_m^{\text{mask}}, \mathbf{S}_m^{\text{mask}})/\tau)}{\sum_{l=1}^M \exp(f(\mathbf{T}_m^{\text{mask}}, \mathbf{S}_l^{\text{mask}})/\tau)}$$

- Masked language model exploits multi-modal fusion.
- Contrastive learning is to learn cross-modal alignments.

# Approach : TAVP



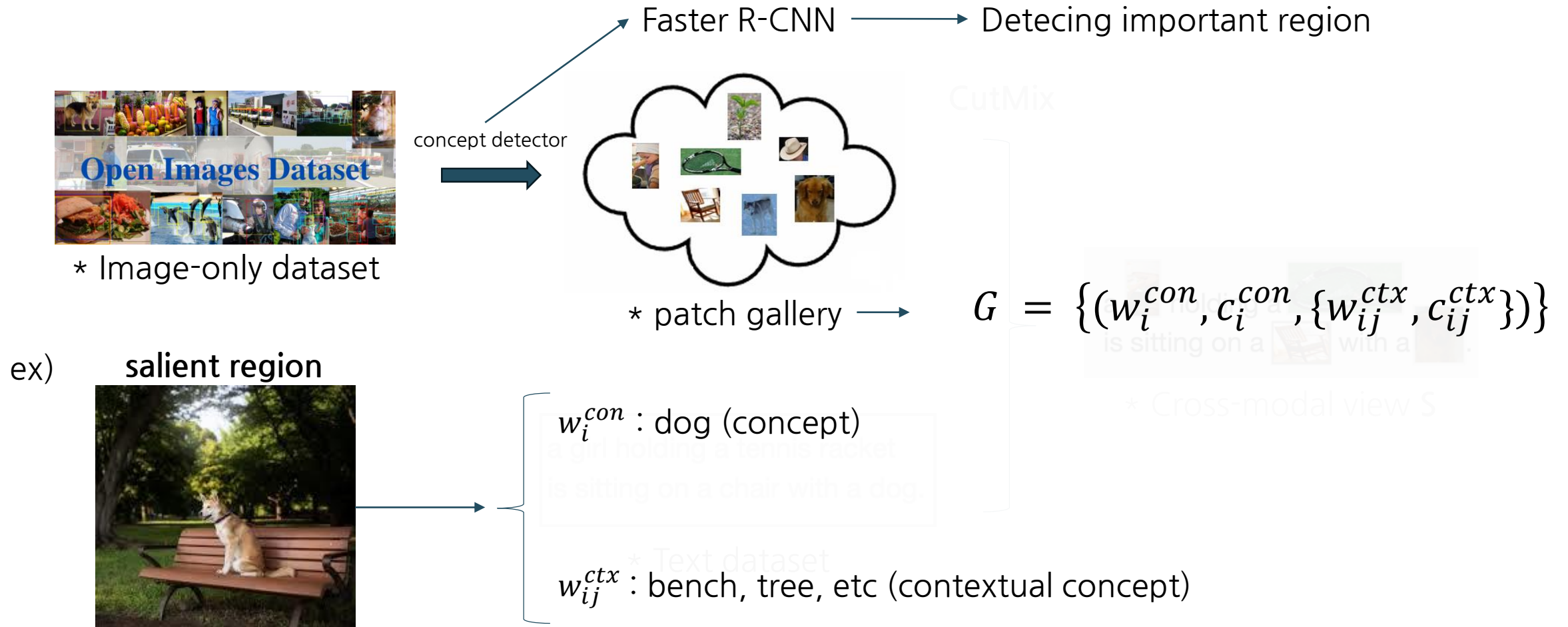
► Effectively combining image and text information using tags from image  $I$

# 3

## Details

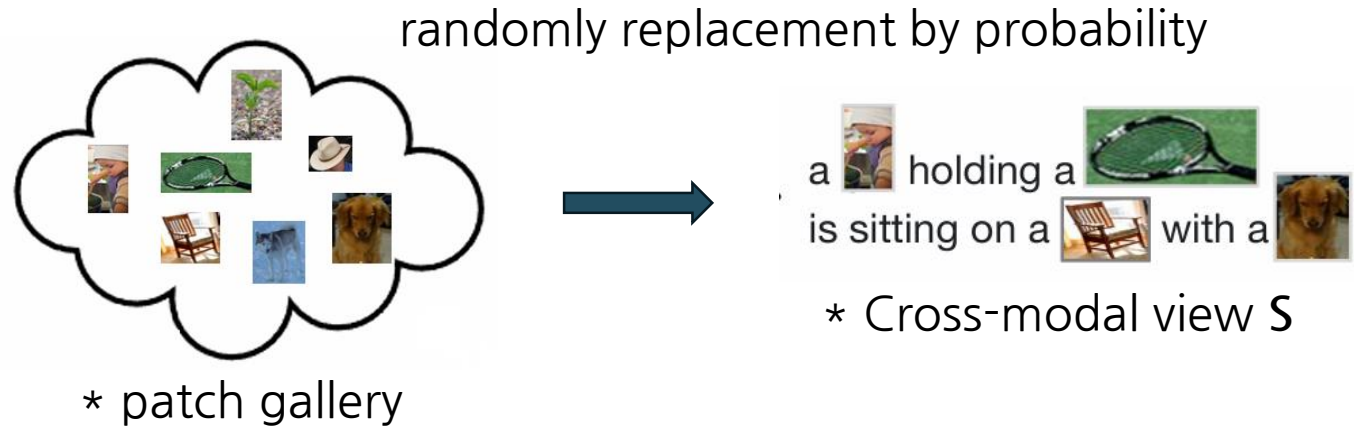
---

# Details : Cross-modal CutMix



► Faster R-CNN : Bottom up attention mechanism

►  $w$  is **concept label** and  $c$  is **confidence score**.



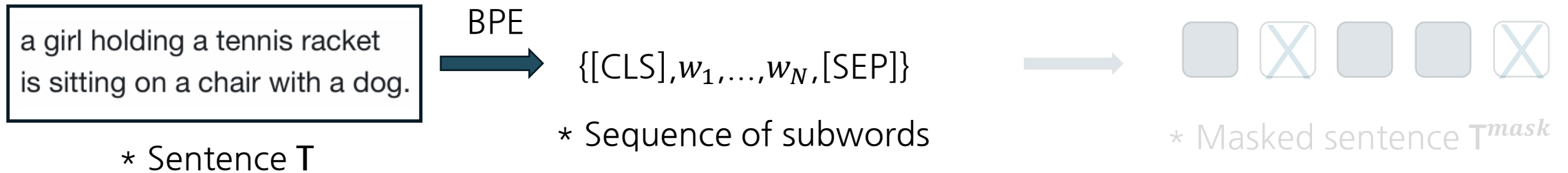
Control the importance of contextual concepts

$$p_i = \begin{cases} c_i^{\text{con}} + \frac{r_{\text{ctx}}}{|G_i|} \sum_{w_{ij}^{\text{ctx}} \in G_i} c_{ij}^{\text{ctx}}, & \text{if } w_i^{\text{con}} = w_n \\ 0, & \text{otherwise} \end{cases}$$

The # of set of contextual concepts related to the word

- The text is replaced with patches considering the contextual situation.

# Details : VALP - Byte pair encoding



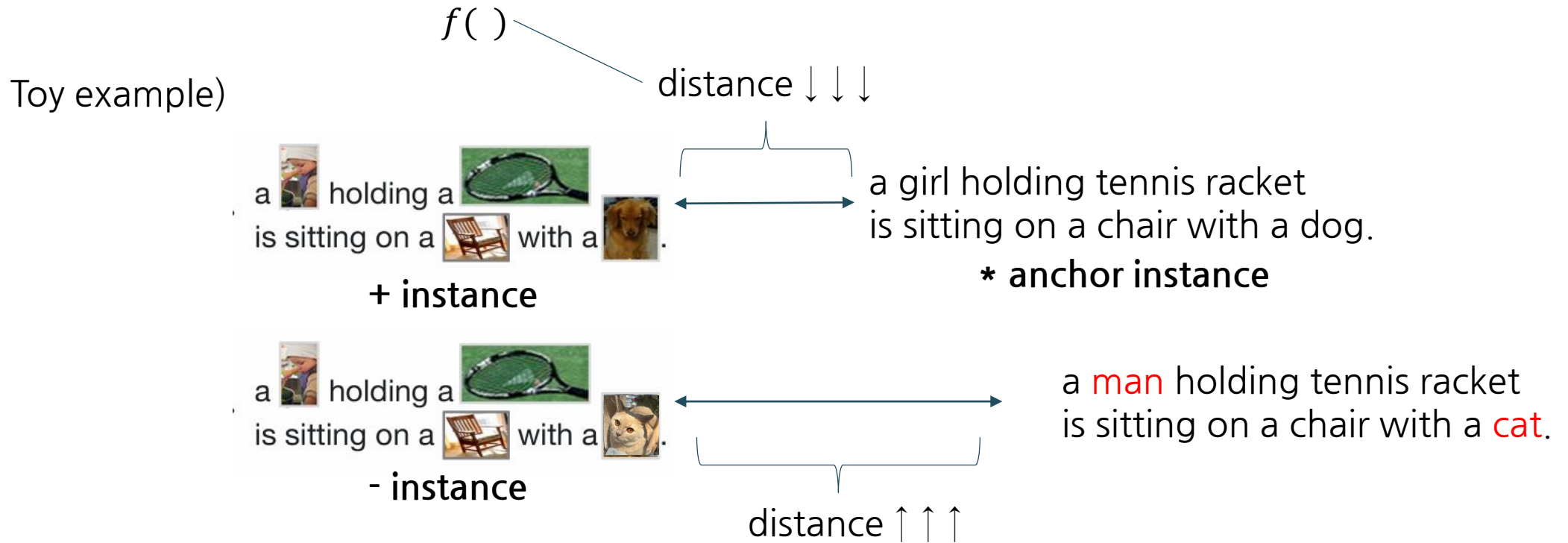
ex)

$\{[CLS], lower, high, [SEP]\} \longrightarrow \{[CLS], low, er, hi, gh, [SEP]\}$

- BPE : word to subword (original  $\rightarrow$  small)
- I think that since patches are smaller than original image, also words will be sliced.

# Details : VALP - Constrative Learning

$$\blacktriangleright \mathcal{L}_{cl} = - \sum_{m=1}^M \log \frac{\exp(f(\mathbf{T}_m^{\text{mask}}, \mathbf{S}_m^{\text{mask}}) / \tau)}{\sum_{l=1}^M \exp(f(\mathbf{T}_m^{\text{mask}}, \mathbf{S}_l^{\text{mask}}) / \tau)}$$



$\blacktriangleright$  I show that not masked example for explanation.



# 4

## Experiment

---

# Datasets

Dataset	Images	Texts	Text Domain
COCO (train)	112K	560K	Image Caption
Conceptual Captions (train)	3M	3M	Image Caption
SBU Caption (all)	840K	840K	Image Caption
Flickr30k (train)	29K	145K	Image Caption
VQA (train)	83K	445K	Question
GQA (train)	79K	1.0M	Question
VG-QA (train)	87K	931K	Question
MSVD (train)	-	48K	Video Caption
MSRVTT (train)	-	130K	Video Caption
VATEX (train)	-	260K	Video Caption
ActivityNet Captions (train)	-	36K	Video Caption
Shutterstock (all)	-	1M	Caption
BookCorpus	-	14M	General Text
OpenImages (od train)	1.67M	-	-

# Setting

## ❖ Backbone Model

- Base Transformer
  - 12 layers of transformer blocks
  - Hidden size :768

## ❖ Position Embedding

- Language/Tag tokens
  - Learning position embedding
- Patch/Image Tokens
  - Linear projection of spatial positions

## ❖ Restriction of token length

- TAVP : 100
- VALP : 80

## ❖ Fine tuning

- VLMixer : BERT - Base
- image and text 300k pre-trained
- lr : 5e-5, mini batch size : 1024
- optim : adam

## ❖ ETC.

- patch detector : Resnet-152 C4
- 15-shot CMC
- $r_{cmc}$  : 0.5
- $r_{ctx}$  : 0.5
- CMCL : temperature ratio : 0.1

# VLMixer: Unpaired Vision-Language Pre-training via Cross-Modal CutMix

Method	Pre-training Data		VQA	NLVR <sup>2</sup>		Text Retrieval			Image Retrieval			GQA
	Image	Text	Test-Dev	Dev	Test	R@1	R@5	R@10	R@1	R@5	R@10	Test-Dev
<b>Paired VLP</b>												
UnicoderVL <sub>base</sub> (Li et al., 2019a)			-	-	-	62.3	87.1	92.8	46.7	76.0	85.3	-
UNITER <sub>base</sub> (Chen et al., 2019)			72.27	77.14	77.87	63.3	87.0	93.1	48.4	76.7	85.9	-
OSCAR <sub>base</sub> (Li et al., 2020b)			73.16	78.07	78.36	70.0	91.1	95.5	54.0	80.8	88.5	61.58
VILT <sub>base</sub> (Kim et al., 2021)			71.26	75.70	76.13	61.5	86.3	92.7	42.7	72.9	83.1	-
VinVL <sub>base</sub> (Zhang et al., 2021)			75.95	82.05	83.08	74.6	92.6	96.3	58.1	83.2	90.1	65.05
ALBEF (Li et al., 2021a)			75.84	82.55	83.14	77.6	94.3	97.2	60.7	84.3	90.5	-
<b>Unpaired VLP</b>												
BERT <sub>base</sub> (Devlin et al., 2019)	None	None	64.85	51.30	51.34	57.44	84.00	91.58	44.03	74.12	84.06	50.20
VinVL <sub>unpaired</sub> (Zhang et al., 2021)	COCO	COCO	71.78	71.14	72.01	61.92	86.90	93.08	46.90	76.18	85.53	62.24
U-VisualBERT (Li et al., 2021b)*	COCO	COCO	72.41	-	-	-	-	-	-	-	-	-
VLMixer	COCO	COCO	<b>72.60</b>	<b>72.71</b>	<b>73.08</b>	<b>62.69</b>	<b>87.35</b>	<b>93.64</b>	<b>47.95</b>	<b>77.06</b>	<b>86.22</b>	<b>63.13</b>
U-VisualBERT (Li et al., 2021b)	CC3M	CC3M+BC	70.74	71.74	71.02	-	-	-	-	-	-	-
VinVL <sub>unpaired</sub> (Zhang et al., 2021)	CC3M	CC3M	72.20	68.96	68.94	62.08	86.04	93.00	47.29	76.15	85.53	63.12
VLMixer	CC3M	CC3M	72.66	74.31	73.86	62.20	86.32	92.80	47.44	76.22	85.41	62.65
VLMixer	Full	Full	<b>72.89</b>	<b>76.61</b>	<b>77.01</b>	<b>64.76</b>	<b>88.56</b>	<b>94.22</b>	<b>50.06</b>	<b>78.36</b>	<b>86.91</b>	<b>63.25</b>