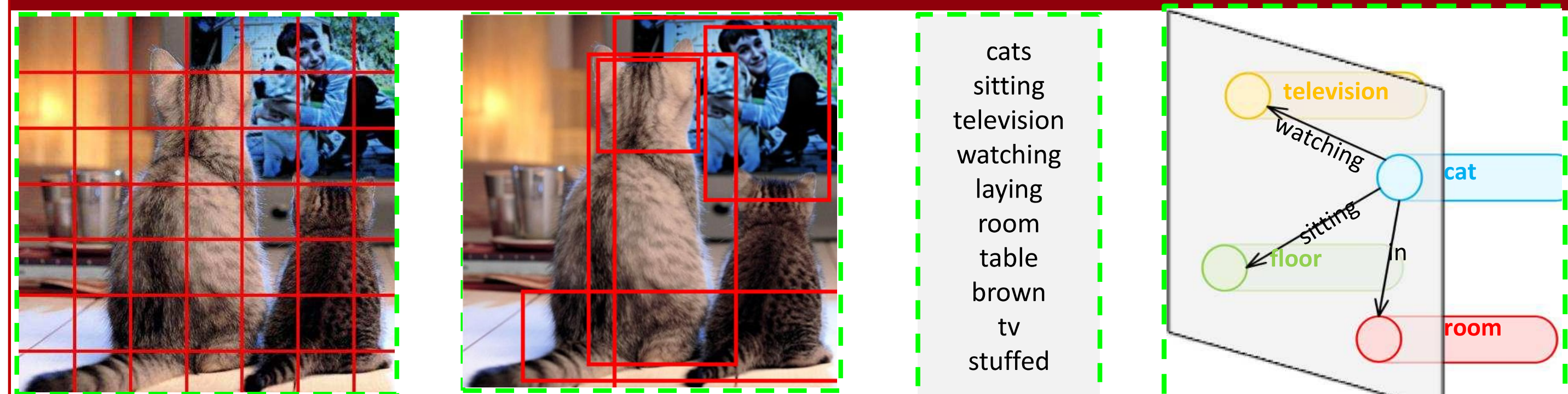




## Introduction



**Figure 1.** Illustrations of commonly-used image representations (from left to right): CNN-based grid visual features, RCNN-based region visual features, textual concepts, and scene-graphs.

An image in vision-and-language tasks, e.g., image captioning and visual question answering, is typically represented in two fundamental forms: **visual features** and **textual concepts** (see Figure 1).

## Limitation & Challenge:

- Most existing downstream systems integrate visual features and the textual concepts **in the decoding process**, mostly ignoring **the innate alignment** between the two modalities.
- The systems have to learn the **alignment** between each individual visual feature and textual concept.
- These representations only contain **individual features**, lacking the **meaningful combinations** and **structural relationships** among them.

Those problems hinder the system from understanding images efficiently.

## Solution:

- We propose the **Mutual Iterative Attention (MIA)** module to **align** the visual features and textual concepts **in the encoding process**. Using textual concepts to **query** and **integrate** visual features with attention, we could get image representations centered upon each concept forming **meaning visual feature groups**, and vice versa. The representations are refined by applying MIA **iteratively**.

## Approach

Our approach based on the Multi-Head Attention (MHA) and Feed-Forward Network (FCN) from Transformer [1].

### Mutual Attention

Given visual features  $I$  and textual concepts  $T$ , the mutual attention is conducted as:

$$I' = \text{FCN}(\text{MHA}(T, I)), \quad T' = \text{FCN}(\text{MHA}(I', T)) \quad (1)$$

i.e., visual features are first **integrated** according to textual concepts, and then textual concepts are **integrated** according to integrated visual features.

### Mutual Iterative Attention (MIA)

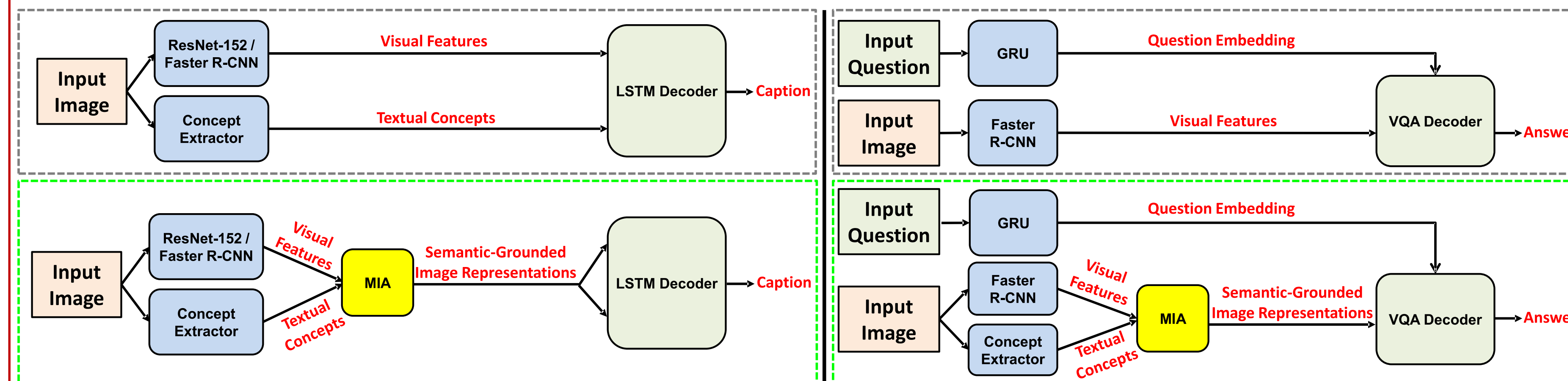
We perform mutual attention **iteratively** to refine both visual features and textual concepts:

$$I_N = \text{FCN}(\text{MHA}(T_{N-1}, I_{N-1})), \quad T_N = \text{FCN}(\text{MHA}(I_N, T_{N-1})) \quad (2)$$

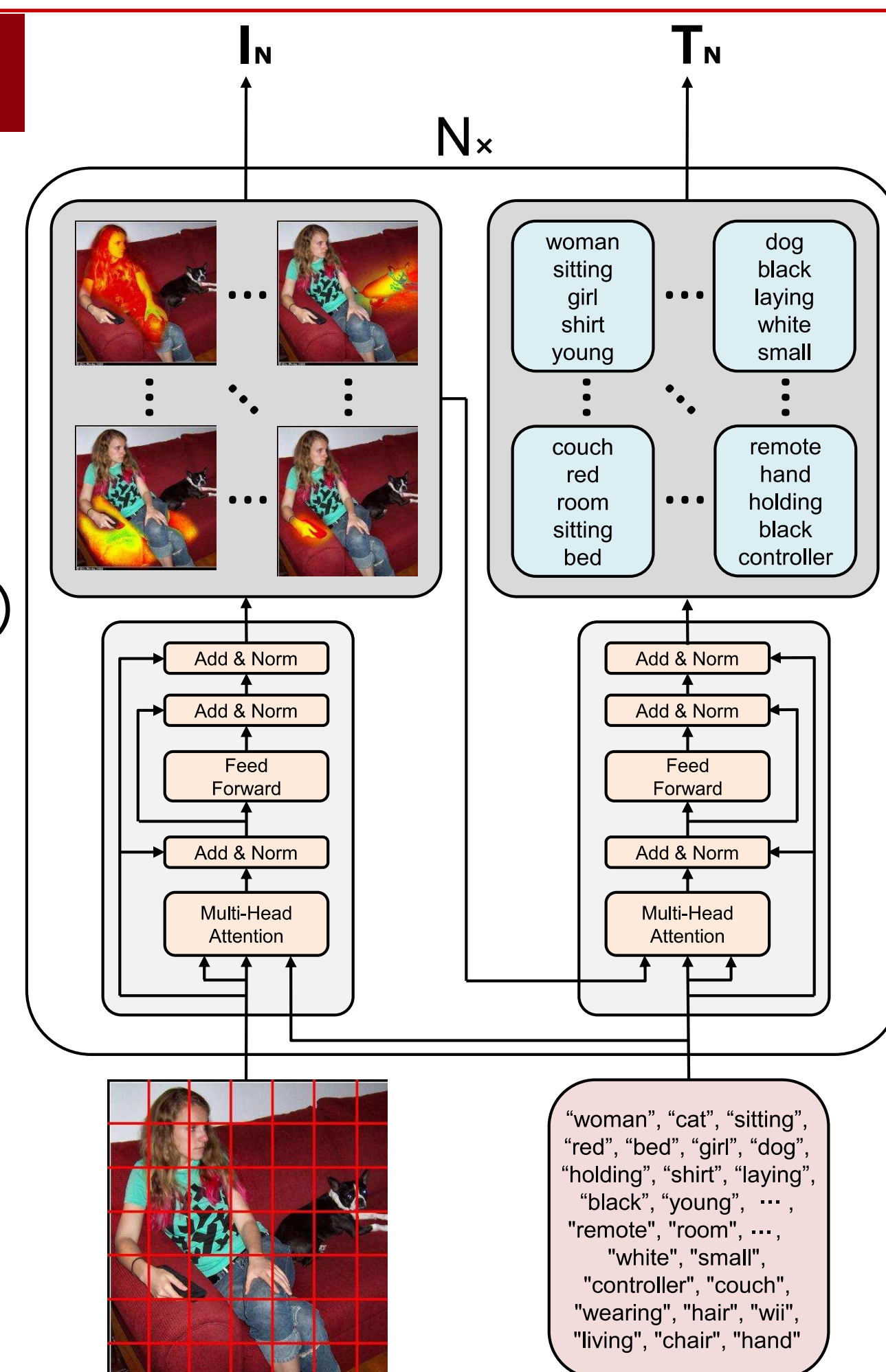
### Semantic-Grounded Image Representations

Since the visual features and the textual concepts are already **aligned**, we can add them up to get the **semantic-grounded image representations**:

$$\text{MIA}(I, T) = \text{LayerNorm}(I_N + T_N) \quad (3)$$



**Figure 3.** Illustration of how to equip the baseline models with our MIA. MIA aligns and integrates the original image representations from two modalities. Left: For image captioning, the semantic-grounded image representations are used to replace both kinds of original image features. Right: For VQA, MIA only substitutes the image representations, and the question representations are preserved.



**Figure 2.** MIA combines the individual features (the lower) from each domain, resulting in integrated image representations (the upper) reflecting certain semantics of the image.

## Experiments

- We evaluate the proposed MIA on two multi-modal tasks (image captioning and visual question answering (VQA)).

Methods	B-1	B-2	B-3	B-4	M	R	C	S
Visual Attention	72.6	56.0	42.2	31.7	26.5	54.6	103.0	19.3
w/ MIA	<b>74.5</b>	<b>58.4</b>	<b>44.4</b>	<b>33.6</b>	<b>26.8</b>	<b>55.8</b>	<b>106.7</b>	<b>20.1</b>
Concept Attention	72.6	55.9	42.5	32.5	26.5	54.4	103.2	19.4
w/ MIA	<b>73.8</b>	<b>57.4</b>	<b>43.8</b>	<b>33.6</b>	<b>27.1</b>	<b>55.3</b>	<b>107.9</b>	<b>20.3</b>
Visual Condition	73.3	56.9	43.4	33.0	26.8	54.8	105.2	19.5
w/ MIA	<b>73.9</b>	<b>57.3</b>	<b>43.9</b>	<b>33.7</b>	<b>26.9</b>	<b>55.1</b>	<b>107.2</b>	<b>19.8</b>
Concept Condition	72.9	56.2	42.8	32.7	26.4	54.4	104.4	19.3
w/ MIA	<b>73.9</b>	<b>57.3</b>	<b>43.9</b>	<b>33.7</b>	<b>26.9</b>	<b>55.1</b>	<b>107.2</b>	<b>19.8</b>
Visual Regional Attention	75.2	58.9	45.2	34.7	27.6	56.0	111.2	20.6
w/ MIA	<b>75.6</b>	<b>59.4</b>	<b>45.7</b>	<b>35.4</b>	<b>28.0</b>	<b>56.4</b>	<b>114.1</b>	<b>21.1</b>

**Table 1.** Results of representative systems on the MSCOCO image captioning dataset. B-n, M, R, C and S are short for BLEU-n, METEOR, ROUGE-L, CIDEr and SPICE, respectively.

Methods	B	M	R	C	S	Methods	Test-dev	Test-std
Up-Down[2]	36.5	28.0	57.0	120.9	21.5	Up-Down[2]	67.3	67.5
w/ MIA	<b>37.0</b>	<b>28.2</b>	<b>57.4</b>	<b>122.2</b>	<b>21.7</b>	w/ MIA	<b>68.8</b>	<b>69.1</b>
Transformer	39.0	28.4	58.6	126.3	21.7	BAN[3]	69.6	69.8
w/ MIA	<b>39.5</b>	<b>29.0</b>	<b>58.7</b>	<b>129.6</b>	<b>22.7</b>	w/ MIA	<b>70.2</b>	<b>70.3</b>

**Table 2.** Results of systems under the reinforcement learning setting. **Table 3.** The overall accuracy on the VQA task.

- As we can see, the proposed MIA exhibits compelling effectiveness in boosting the baseline systems.

## References

- [1] Attention is all you need. *In NIPS*, 2017.
- [2] Bottom-up and top-down attention for image captioning and VQA. *In CVPR*, 2018.
- [3] Bilinear attention networks. *In NeurIPS*, 2018.

## Contact Us

- fenglinliu98@pku.edu.cn
- liuyuanxin@iie.ac.cn
- renxc@pku.edu.cn,
- xusun@pku.edu.cn
- xiaodong.he@jd.com



➤ arxiv



➤ code