

LQMFormer: Language-aware Query Mask Transformer for Referring Image Segmentation

Supplementary Material

001 1. Gaussian Enhanced Multi-Modal Fusion 002 Module (GMMF)

003 The application of the discrete Fast Fourier Transform (FFT)
004 within the GMMF module is central to transforming the
005 vision-language features F'_{vl} into the frequency domain, for
006 efficient global modelling. The FFT is defined as:

$$007 F_{\text{FFT}}(F'_{vl}) = \mathcal{F}(F'_{vl}) \quad (1)$$

008 where \mathcal{F} denotes the Fourier transform operation, and
009 F'_{vl} are the vision-language features that are obtained from
010 vision backbone model.

011 Amplitude and Phase Components in Fourier 012 Transform

013 The complex exponential in $\mathcal{F}(F'_{vl})$ contains the phase in-
014 formation. The amplitude $\mathcal{A}(F'_{vl})$ and phase $\mathcal{P}(F'_{vl})$ compo-
015 nents, critical for reconstructing and manipulating the vision-
016 language features in the frequency domain, are derived from
017 $\mathcal{F}(F'_{vl})$ and defined as follows:

$$018 \mathcal{A}(F'_{vl}) = \sqrt{\mathcal{R}(F'_{vl})^2 + \mathcal{I}(F'_{vl})^2}, \quad (2)$$

$$019 \mathcal{P}(F'_{vl}) = \arctan \left[\frac{\mathcal{I}(F'_{vl})}{\mathcal{R}(F'_{vl})} \right], \quad (3)$$

021 where $\mathcal{R}(F'_{vl})$ and $\mathcal{I}(F'_{vl})$ represent the real and imaginary
022 components of $\mathcal{F}(F'_{vl})$, respectively.

023 Amplitude Modulation

024 In the GMMF module, the amplitude component $\mathcal{A}(F'_{vl})$
025 is modulated to enhance global visual-language features,
026 specifically using Gaussian smoothed filters for low-pass
027 filtering [1, 3]:

$$028 \mathcal{A}'(F'_{vl}) = \mathcal{A}(F'_{vl}) * \phi(F'_{vl}, \beta), \quad (4)$$

029 where $*$ represents the operation of low-pass filtering, and
030 the result $\mathcal{A}'(F'_{vl})$ is the low-pass filtered version of the am-
031 plitude component and β is learnable bandwidth parameter,
032 calculated from F_{vl} followed by sequence of Linear and
033 Pooling operation.

034 Now, for reconstructing the complex frequency represen-
035 tation from the amplitude and phase, the following equation
036 is used:

$$037 \mathcal{F}'(F'_{vl}) = \mathcal{A}'(F'_{vl}) * e^{j\mathcal{P}(F'_{vl})}, \quad (5)$$

where $e^{j\mathcal{P}(F'_{vl})}$ represents the complex exponential with the
phase component $\mathcal{P}(F'_{vl})$. 038 039

To reconstruct the enhanced features for further process-
ing, following [3], we apply the inverse discrete Fast Fourier
transform (Inverse FFT): 040 041 042

$$F_{vl} = F_{\text{FFT}}^{-1}(\text{Conv}(\mathcal{F}'(F'_{vl}))) + F'_{vl}, \quad (6) \quad 043$$

where F_{FFT}^{-1} denotes the inverse Fourier transform opera-
tion, and Conv indicates 1x1 convolution. The final step in
the GMMF module involves the combination of the original
features F'_{vl} and the Gaussian-enhanced features, here rep-
resented by F_{vl} , resulting from the convolution and inverse
FFT of the modulated complex frequency representation
 $\mathcal{F}'(F'_{vl})$. 044 045 046 047 048 049 050

This reconstruction is critical in the GMMF module, par-
ticularly for modulating the amplitude component while
preserving the phase information, leading to enhanced visual-
language feature representation. This detailed section suppl-
ements method section by offering a more detailed explana-
tion of the operations within the GMMF module, particularly
focusing on the Fourier transform applications and the ratio-
nale behind the use of Gaussian smoothing in the frequency
domain for feature enhancement. 051 052 053 054 055 056 057 058 059

060 2. Qualitative Results

The qualitative analysis shows the LQMFormer capabilities
in segmenting objects from complex visual scenes based on
intricate language descriptions. Each row in Figure 1 reflects
a different scenario where the model proficiency in vision
grounding and language understanding is evaluated. 061 062 063 064 065

Contextual Differentiation: The first row illustrates the
model’s ability to distinguish between individuals based on
their roles and attire – a critical skill in scenes with multiple
similar subjects. Notably, it segments the batter in white and
the referee in black, showing understanding of context based
description. 066 067 068 069 070 071

Detailed Descriptive Segmentation: The second row
presents a case where LQMFormer precisely segments a
batter based on both a descriptive action (“with one knee
on the ground”) and a specific object relation (“a bat on the
ground in front of him”). This shows the model’s capacity
to interpret complex activities and relative positioning. 072 073 074 075 076 077

Disambiguation of Multiple Entities: In the third row,
LQMFormer adeptly segments both a man and a woman
wearing purple, accurately understanding multiple entities 078 079 080



Figure 1. Qualitative Comparison of our model LQMFormer with ReLA [2] on GRES dataset.

081 based on their clothes color, which is a common challenge
082 in crowded scenes.

083 **Understanding of Complex Indirect Descriptions:** The
084 fourth row demonstrates LQMFormer’s ability to under-
085 stand indirect descriptors like “crouching dude” in the fore-
086 ground, effectively linking the alias to the complex descrip-
087 tion in the scene.

088 **Challenging Case of Recognition of Occluded Objects:**
089 Conversely, the fifth row presents a challenging scenario
090 for LQMFormer. The model encounters difficulty accu-
091 rately segmenting an occluded person characterized solely
092 by a color descriptor (“pink sweater”). This case shows
093 the model’s current limitations in distinguishing occluded
094 objects, especially when they blend into the background,
095 indicating a potential direction for future enhancements.

096 These qualitative evaluations shows that LQMFormer is
097 proficient in scenarios that require detailed language under-
098 standing and vision grounding. The model understands a
099 range of expressions, from straightforward references to a
100 single entity to complex descriptions involving multiple ob-
101 jects. Its performance passes compared to previous methods
102 in RIS, highlighting its capability in this complex task.

103 References

- 104 [1] Chongyi Li, Chun-Le Guo, Man Zhou, Zhexin Liang,
105 Shangchen Zhou, Ruicheng Feng, and Chen Change Loy. Em-
106 bedding fourier for ultra-high-definition low-light image en-
107 hancement. *arXiv preprint arXiv:2302.11831*, 2023. 1
- 108 [2] Chang Liu, Henghui Ding, and Xudong Jiang. Gres: Gener-
109 alized referring expression segmentation. In *Proceedings of*
110 *the IEEE/CVF Conference on Computer Vision and Pattern*
111 *Recognition*, pages 23592–23601, 2023. 2
- 112 [3] Bo Miao, Mohammed Bennamoun, Yongsheng Gao, and Aj-
113 mal Mian. Spectrum-guided multi-granularity referring video
114 object segmentation. In *Proceedings of the IEEE/CVF Interna-*
115 *tional Conference on Computer Vision*, pages 920–930, 2023.
116 1