

Multimodal Output Summary with relevant Image

Image-to-Summary

M3LS Dataset

Text-Vision Fusion

Embeddings

Text

ResNet50

Images

Positional  
Encoding

MA

MA

FFNN

$\psi_T^l$

FFNN

$\psi_V^h$

MM Learning

Target  
Summary in  
Input Language

Multilingual  
Decoder

MS Generation

