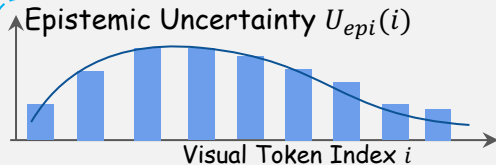
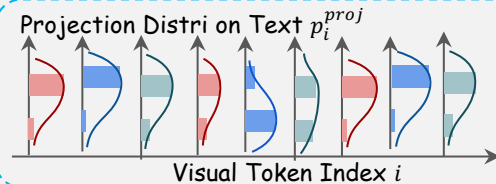


Before Decoding



Uncertainty
Quantification



Vision-to-text
projection

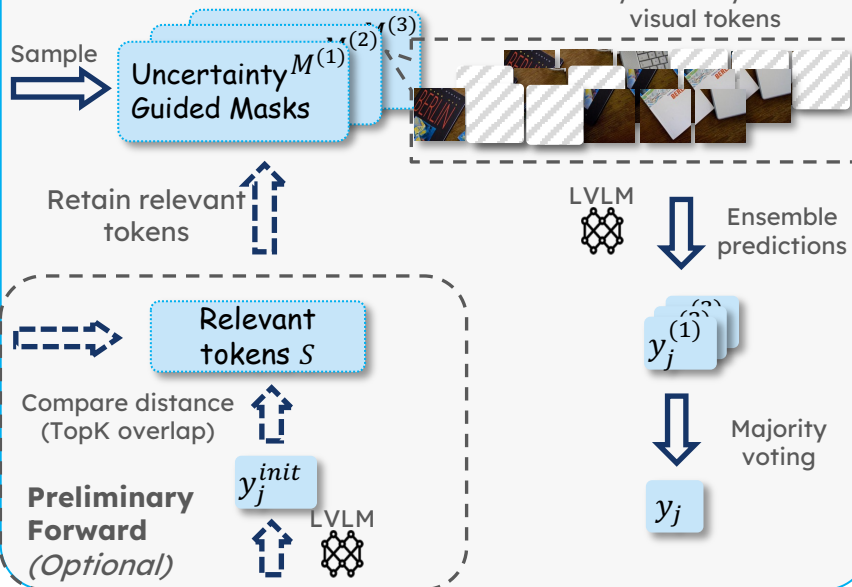
LM decoder
in LVLM

Forward with only visual input



Tokenized Visual Inputs To LM Decoder

Decoding Time at each token j



User: Describe the image. Assistant: y_1, y_2, \dots, y_{j-1}

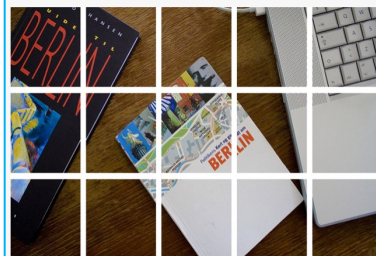
Text Inputs & Previous Generation

...In addition to the book,
there is a **laptop**
positioned to the right of
the book ...

Our SURE-Decoding
(No Hallucination)

...In addition to the book,
there is a laptop **and a**
mouse on the desk...

Original Generation
(Hallucination on "mouse")



Visual Inputs