(a) LLaMA Vision Architecture

(b) Default Attention Masks

(c) Attention Masking Mechanism to analyze different information flows