logit$(x_t | x^v, x^a, x^l)$ → **Entropy-guided Adaptive Decoding** → AVCD OFF Prediction: "Three"

Nine
Six
Eleven
Seven

**Modality Dominance**

$A_{Q,\text{video}}$ $A_{Q,\text{audio}}$ $A_{Q,\text{Language}}$

$Eq. (3)$

**Language** > Video > Audio

**Question**
"How many types of musical instruments sound in the video?"

**Attentive Masking**

Video Masking — BLIND — logit$(x_t | x^{\neg v}, x^a, x^l)$

Audio Masking — logit$(x_t | x^v, x^{\neg a}, x^l)$

AV Masking — BLIND — logit$(x_t | x^{\neg v}, x^{\neg a}, x^l)$

AVCD On Prediction: **"Two"**

Three
One
Two

$Eq. (10)$

Three
One
Two

Ground Truth: "Two"