AVVS Previous Method (i) Seperate-Dependence Fusion Input 1: Video video's temporal dependence Input 2: Audio Audio-visual interaction Frame 1~3: Frame 4: Frame 5: Violin / Singing Singing Piano/Singing (ii) Object-Limited Queryless Decoding Frame 4: Mask of person Output: (Sound of Violin and Singing) **FCN** Decoder Frame 1~3: Frame 5:Mask of Mask of violin and person piano and person Piano segmented by mistake (b) (a)

