| Audio | GT | Pose video $V_P$ | Face&hand mask $V_M^{f+h}$ | Lips mask $V_M^l$ | Backgroud mask $V_M^b$ |