

Audiovisual speech recognition (AV-ASR) utilizes the movement of the speaker's lips and mouth region as compensation for acoustic information to recognize the spoken utterances. Compared to pure acoustic ASR, AV-ASR has shown its superior performance in acoustically noisy environments. Advancements in English AV-ASR have been largely driven by the availability of large public English audiovisual datasets. Conversely, other languages such as German are currently considered low-resource for AV-ASR tasks, lacking sizable labeled audiovisual datasets. To address this challenge, we employed transfer learning for the German AV-ASR task, where the self-supervised model pre-trained on unlabeled English data is fine-tuned on German data. Our results from the SLT 2022 Hackathon demonstrate the efficacy of transfer learning for low-resource AV-ASR tasks. Further research will focus on the effective transfer learning approaches and the interpretability of these transfer learning techniques.