

| Task | Model        | Metric                           | Public Test Set                                  | Public Result                | Internal Test Set    | Internal Result |
|------|--------------|----------------------------------|--|------------------------------|----------------------|-----------------|
| SRWT | Whisper-L-v3 | AAS(ms, ↓)<br>(Shi et al., 2022) | The comparison plan test set is not open source. | —                            | Test <sub>srwt</sub> | 11.09           |
|      | Qwen-Audio   |                                  |  | —                            |                      | 9.17            |
|      | GMM-HMM      |                                  |  | —                            |                      | <b>7.55</b>     |
|      | OSUM         |                                  |  | —                            |                      | 7.61            |
| VED  | Qwen2-Audio  | ACC<br>(%, ↑)                    | ESC-50<br>VocalSound                             | -<br>93.3                    | Test <sub>ved</sub>  | 33.25           |
|      | TouchASP     |                                  |  | 85.7<br>-                    |                      | —               |
|      | PANNs        |                                  |  | 83.3<br>-                    |                      | 3.25            |
|      | OSUM         |                                  |  | 96.02<br>80.75               |                      | <b>81.00</b>    |
|      | Qwen2-Audio  |                                  |  | 55.3<br>-                    |                      | 38.04           |
|      | TouchASP     |                                  |  | 50.5<br>-                    |                      | —               |
| SER  | Sensevoice-L | ACC<br>(%, ↑)                    | MELD<br>test<br>MER2023<br>test                  | <b>63.1</b><br>69.2          | Test <sub>ser</sub>  | —               |
|      | Sensevoice-S |                                  |  | 57.8<br>68.3                 |                      | 40.77           |
|      | Emotion2Vec  |                                  |  | 51.88<br>—                   |                      | 51.14           |
|      | EmoBox       |                                  |  | 51.89<br>65.23               |                      | <b>74.54</b>    |
|      | OSUM         |                                  |  | 56.64<br><b>88.84</b>        |                      | 72.97           |
|      | GLM-4        |                                  |  | —                            |                      | 53.97           |
|      | OSUM         |                                  |  | —                            |                      | <b>68.56</b>    |
| SGC  | Qwen2-Audio  | ACC<br>(%, ↑)                    | AISHELL-1<br>test                                | <b>97.36</b><br><b>97.25</b> | Test <sub>sgc</sub>  | <b>98.43</b>    |
|      | OSUM         |                                  | Kaggle-CommonVoice<br>valid-test                 | <b>100</b><br><b>98.75</b>   |                      | 93.74           |
| SAP  | Qwen2-Audio  | ACC<br>(%, ↑)                    | Kaggle-CommonVoice<br>valid-test                 | <b>35.53</b>                 | Test <sub>sap</sub>  | 49.52           |
|      | OSUM         |                                  |  | <b>76.72</b>                 |                      | <b>67.94</b>    |
| STTC | Qwen2-Audio  | GPT-3.5-Turbo<br>Scoring         | AirBench<br>speech                               | <b>6.77</b>                  | Test <sub>sttc</sub> | <b>6.91</b>     |
|      | ASLP-Audio   |                                  |  | 4.96                         |                      | 6.69            |