

| Question   | Hypothesis / Factor         | Option   | %  |
|--|-----------------------------|--|----|
| Please listen to the 3 audios below completely<br><br>1. Reference audio<br>2. System A<br>3. System B<br><br>Kindly mark your preference.   | Reference-Matching          | System B is better than System A because System B is closer to the Reference                                 | 55 |
|  |                             | System A is better than System B because System A is closer to the Reference                                 | 20 |
|  |                             | System B is better than System A in general.   | 11 |
|  |                             | System A is better than System B in general.   | 7  |
|  |                             | Can 't find the difference.  | 7  |
| Please listen to the 3 audios below completely<br><br>1. Reference audio<br>2. System A<br>3. System B<br><br>Compare System A and System B against with the reference. System A made a pronunciation mistake at 2 words but had excellent voice quality. System B had average digital voice quality but showed excellent pronunciation with no mistakes. Which option would you choose? | Judgement Ambiguity         | I would rate System A better than System B.  | 56 |
|  |                             | I would rate System B better than System A.  | 33 |
|  |                             | I would rate System A equal to System B.   | 11 |
| How often did you come across confusing situations as the above?   | Confounders                 | Only Few times did I face difficulty in assigning scores to different systems.                               | 56 |
|  |                             | Always clear how to rate the different systems.  | 31 |
|  |                             | Many times faced difficulty in assigning scores to different systems due to ambiguity in scoring guidelines. | 7  |
|  |                             | Every page I faced difficulty in assigning scores to different systems.                                      | 6  |
| Please listen to the 3 reference audios below completely<br><br>1. Reference 1<br>2. Reference 2<br>3. Reference 3<br><br>On average, I felt the reference audio samples I listened to were -  | Label Ambiguity [Reference] | Excellent (100 - 80)   | 42 |
|  |                             | Good (80 - 60)   | 38 |
|  |                             | Fair (60 -40)  | 19 |
|  |                             | Poor (40 - 20)   | 1  |
| How would you rate below sample generated System A?<br><br>1. System A   | Label Ambiguity [System]    | Good (60 - 80)   | 45 |
|  |                             | Excellent (80 -100)  | 35 |
|  |                             | Fair (40 -60)  | 12 |
|  |                             | Poor (20 - 40)   | 6  |
|  |                             | Bad (0 - 20)   | 2  |
| On a scale of 5, how difficult was the MUSHRA test?  | Difficulty                  | Moderate   | 40 |
|  |                             | Easy   | 34 |
|  |                             | Difficult  | 12 |
|  |                             | Very easy  | 10 |
|  |                             | Very difficult.  | 3  |