

The following publication Yan, J., Chen, F., Gao, X., & Peng, G. (2021). Auditory-Motor Mapping Training Facilitates Speech and Word Learning in Tone Language-Speaking Children With Autism: An Early Efficacy Study. *Journal of Speech, Language, and Hearing Research*, 64(12), 4664-4681 is available at https://dx.doi.org/10.1044/2021_JSLHR-21-00029. The *Journal of Speech, Language, and Hearing Research* is available at <https://pubs.asha.org/journal/jslhr>.

1 **Auditory-Motor Mapping Training Facilitates Speech and Word Learning in Tone-**
2 **Language-Speaking Children with Autism: An Early Efficacy Study**

3 Jinting Yan^{1#}, Fei Chen^{2#, *}, Xiaotian Gao¹, and Gang Peng³

4 ¹ *College of Qiyue Communication & Cangzhou Research Centre for Child Language*
5 *Rehabilitation, Cangzhou Normal University, Cangzhou, China*

6 ² *School of Foreign Languages, Hunan University, Changsha, China;*

7 ³ *Research Centre for Language, Cognition, and Neuroscience & Department of Chinese and*
8 *Bilingual Studies, The Hong Kong Polytechnic University, Hong Kong SAR*

9 **#The first two authors contribute equally to this study.**

10 *** Corresponding author: Fei Chen**

11 School of Foreign Languages, Hunan University, Changsha, China.

12 **E-mail addresses:** chenfeianthony@gmail.com

13

14 **Conflict of Interest:** The authors have declared that no competing interests existed at the time of
15 publication.

16 **Funding Statement:** This work was in part supported by the Major Program of National Social
17 Science Foundation of China (18ZDA293), the Natural Science Foundation of China overseas
18 collaboration grant (31728009), the General Research Fund of the Research Grants Council of
19 Hong Kong (15610321), and a key project from *Cangzhou Normal University* (xnjjw1907).

20 **Abstract**

21 **Purpose:** It has been reported that tone-language-speaking children with autism demonstrate speech-
22 specific lexical tone processing difficulty although they have intact or even better-than-normal processing
23 of nonspeech/melodic pitch analogues. In this early efficacy study, we evaluated the therapeutic potential
24 of an Auditory-Motor Mapping Training (AMMT) in facilitating speech and word output for Mandarin-
25 speaking nonverbal and low-verbal children with autism, in comparison with a matched non-AMMT-based
26 control treatment.

27 **Method:** Fifteen Mandarin-speaking nonverbal and low-verbal ASD children participated and completed
28 all the AMMT-based treatment sessions by intoning (singing) and tapping the target words delivered via
29 an app, while another fifteen participants received control treatment. Generalized linear mixed-effects
30 models were created to evaluate the speech production accuracy and word production intelligibility across
31 different groups and conditions.

32 **Results:** Results showed that the AMMT-based treatment provided a more effective training approach in
33 accelerating the rate of speech (especially lexical tone) and word learning in the trained items. More
34 importantly, the enhanced training efficacy on lexical tone acquisition remained at two weeks post-therapy,
35 and generalized to untrained tones that were not practiced. Furthermore, the low-verbal participants showed
36 higher improvement compared to the nonverbal participants.

37 **Conclusions:** These data provide the first empirical evidence for adopting the AMMT-based training to
38 facilitate speech and word learning in Mandarin-speaking nonverbal and low-verbal children with autism.
39 This early efficacy study holds promise for improving lexical tone production in Mandarin-speaking
40 children with autism but should be further replicated in larger-scale randomized studies.

41 **Keywords:** autism, nonverbal, low-verbal, Mandarin, lexical tones, Melodic Intonation Therapy

42 **Auditory-Motor Mapping Training Facilitates Speech and Word Learning in Tone-**
43 **Language-Speaking Children with Autism: An Early Efficacy Study**

44
45 **1 Introduction**

46 Autism spectrum disorder (ASD) is a neurodevelopmental disorder identified by a
47 constellation of early-appearing social communication deficits and restricted, repetitive sensory-
48 motor behaviours (American Psychiatric Association, 2013). The delayed/atypical language
49 development – historically linked to a defining feature – has been removed from the diagnostic
50 criteria and is now regarded as one of the co-occurring characteristics of autism. The goal of this
51 study is to integrate recent advances in speech therapy to facilitate speech and word learning in
52 tone-language-speaking children with ASD.

53 *Nonverbal and Minimally Verbal Children with ASD and Behavioral Interventions*

54 Estimates vary, with 25%–46% of children with ASD reported as minimally verbal past
55 the age of 5 years (Norrelgen et al., 2015; Rose et al., 2016; Tager-Flusberg & Kasari, 2013),
56 meaning that they have an expressive vocabulary fewer than 20 (Kasari et al., 2013), or 30
57 intelligible words (Bal et al., 2016; Plesa-Skwerer et al., 2016). Besides, some individuals with
58 ASD even remain nonverbal with a complete absence of functional speech, and lack the ability to
59 communicate with others using spoken language (Klinger et al., 2002; Koegel et al., 2009;
60 Turner et al., 2006). While children with ASD showed deficits in social communication and
61 social interaction, this is further exacerbated in nonverbal and minimally verbal children with
62 ASD due to their severely delayed speech development. Increased speech and word output is
63 considered a positive prognostic indicator of outcomes for nonverbal or minimally verbal
64 children with ASD (Lord et al., 2006). Recently, the growing availability of computers and

65 smartphones/tablets has generated a great deal of enthusiasm for their potential to help
66 ameliorate speech deficits in minimally-verbal children with ASD, such as the computer-assisted
67 3-D virtual pronunciation tutor (Chen et al., 2019), and the Proloquo2Go iPad app (King et al.,
68 2014). In consideration of autistic features especially among nonverbal and low-verbal
69 individuals who begin treatment with limited or no spoken words, available behavioral
70 interventions often tried to present speech sounds with computerized “high-tech” solutions
71 (Tager-Flusberg & Kasari, 2013).

72 *Mandarin Chinese as a Tonal Language and the State of the ASD Research in China*

73 A very recent review estimated that ASD prevalence in China was comparable to the
74 Western world (about 108 per 10,000) using standardized case identification protocol (Sun et al.,
75 2019). However, the situation of ASD in China lags considerably behind those in the West in
76 terms of public awareness, education opportunities, and life outcomes of people with autism (Liu
77 et al., 2016; Yu et al., 2020). It is compelling to come up with a language-specific training
78 approach for children with ASD who live in a country hosting nearly 20% of the world’s
79 population. Modern Mandarin is a tone language with a relatively simple syllable structure,
80 consisting of initials (consonants), finals (including monophthong, diphthong, triphthong, or
81 nasal finals with a vowel nucleus), as well as lexical tones (Chen et al., 2017). Specifically,
82 Mandarin Chinese is the official and widely spoken language in China and is also used in some
83 other countries/regions, which phonologically differs from English. For example, Mandarin is a
84 syllable-timed tone language that exploits variations in both pitch height and pitch direction at
85 the syllable level to distinguish lexical meanings (Figure 1), while English is a stress-timed non-
86 tone language (Mok & Dellwo, 2008). There are four citation tones in Mandarin Chinese: high-
87 level Tone 1 (T1, [55]), mid-rising Tone 2 (T2, [35]), low-falling-rising Tone 3 (T3, [214]), and

88 high-falling Tone 4 (T4, [51]). These four lexical tones are essential elements of Mandarin
89 speech sounds, and are used to differentiate lexical meanings. For instance, “i” spoken with the
90 four distinct tones can respectively mean “doctor” (T1), “move” (T2), “rely on” (T3), and “easy”
91 (T4). Thus, changing the lexical tone in a tone language has the same kind of effect as changing
92 a vowel or a consonant. In contrast, variation in pitch contours in non-tone languages mainly
93 conveys different moods or intonations without changing the word content (Wang, 1973). Thus,
94 speech therapy in tone-language-speaking children with ASD should not only aim at enhancing
95 consonant and vowel production (segmental elements of speech), but additionally, also aim at
96 improving lexical tone production (supra-segmental element of speech).

97

98 **[Figure 1 about here]**

99

100 *Superior Music and Nonspeech Processing but Impaired Lexical Tone Processing in ASD*

101 Music is one of the most meaningful and popular forms of nonspeech sound; like speech,
102 it has developed to take advantage of the efficiencies of the human auditory system (Baldwin,
103 2012). In the research into auditory processing, several studies demonstrated that children with
104 ASD prefer musical or nonspeech stimuli over speech and attend to them more (Dawson et al.,
105 1998; Kuhl et al., 2005). Many children with ASD showed enhanced capacity of music processing
106 and exhibited strong interest in learning and making music (Buday, 1995; Hairston, 1990; Heaton
107 et al., 1998). Irrespective of the tone inventory (Mandarin or Cantonese), the enhanced or at least
108 preserved non-linguistic pitch perception skill in ASD generalized to those from a tone language
109 background (Cheng et al., 2017; Wang et al., 2017; Yu et al., 2015). Nevertheless, such an
110 enhanced acoustic pitch perception skill was conversely changed to the compromised perception

111 of lexical tones when the pitch changes are phonologically relevant for the tone language speakers
112 with ASD (Chen et al., 2016; Lau et al., 2020; Wang et al., 2017; Yu et al., 2015). In terms of
113 speech production, when compared with age-matched typically-developing (TD) children, the
114 ASD group aged 3–6 years showed an apparent speech delay in the production of Mandarin initials,
115 finals, as well as lexical tones (Wu et al., 2020). Since both music notes and lexical tones share the
116 same psychoacoustical attribute of pitch information, it would be reasonable to take advantage of
117 the relative strength of musical skills to compensate for the relative weakness of speech sound,
118 especially lexical tone acquisition, for tone-language-speaking children with ASD.

119 *The Intonation-Based Interventions in ASD*

120 Given the behavioral resemblance between singing and speaking as well as neural
121 overlap in responses to speech and musical stimuli (Peretz et al., 2015), researchers have begun
122 to examine the therapeutic effects of singing/intonation, and how it can potentially ameliorate
123 some of the speech deficits associated with neurological disorder such as ASD (Wan et al.,
124 2010). Two earlier case studies have described the positive role of singing (intoned rather than
125 spoken verbal stimulus) in facilitating speech development of children with autism (Hoelzley,
126 1993; Miller & Toca, 1979). Recently, Melodic Intonation Therapy (MIT; Albert et al., 1973;
127 Sparks et al., 1974), initially designed for improving spoken language in left-hemisphere stroke
128 patients with severe nonfluent aphasia, has been introduced into speech therapy in ASD. The
129 MIT approach involves the musical elements of both melody and rhythm through the use of
130 pitched vocalization or singing in combination with left-hand rhythmic tapping to provide cueing
131 for syllable production (Norton et al., 2009). An modified version of MIT, called Auditory-
132 Motor Mapping Training (AMMT), was initially proposed by Wan et al. (2011) to facilitate
133 speech output for English-speaking nonverbal children with autism. The intonation-based

134 AMMT combines intonation (singing) of bi-syllabic words or phrases and the use of a pair of
135 tuned drums to activate bimanual motor activities. All six English-speaking nonverbal children
136 with autism who participated in Wan et al.'s (2011) study showed noticeable improvements in
137 their ability to articulate several word approximations. Furthermore, the efficacy of AMMT has
138 been further corroborated in English-speaking minimally verbal children with autism
139 (Chenausky et al., 2016) and one more-verbal child with autism (Chenausky et al., 2017) when
140 compared with a control treatment. In this study, we apply and assess the effect of intonation-
141 based AMMT, with proven efficacy as an intervention for English-speaking children with ASD
142 (Chenausky et al., 2016, 2017; Sandiford et al., 2013; Wan et al., 2011), on Mandarin-speaking
143 children with ASD.

144 To this end, the present study evaluated the therapeutic potential of AMMT in facilitating
145 speech output for tone-language-speaking children with ASD. Moreover, as suggested by Wang
146 (1978), children do not learn speech by acquiring units like phonemes or allophones, but rather by
147 gradually adding lexical items to their repertoire. Two important studies (Ferguson & Farwell,
148 1975; Hsieh, 1972) investigated the development of phonological production in relation to the
149 acquisition of words, and showed that there was a primacy of word learning during speech
150 development. Thus, in this study, our intervention based on a smartphone/iPad app called Music-
151 Mediated and Lexicon-Integrated (MMLI) training, tries to combine speech and word learning.
152 The current training study aimed to assess the efficacy of AMMT-based MMLI in comparison to
153 a matched non-AMMT-based control treatment, Speech Repetition Therapy (SRT) (Chenausky et
154 al., 2016). Specifically, this study aimed to address the following research questions:

- 155 a) Over the intensive treatment sessions of MMLI, would the nonverbal and low-verbal
156 Mandarin-speaking children with ASD show any improvement in speech and word
157 productions?
- 158 b) Compared with the control group, would the MMLI group show a greater improvement in
159 the speech production accuracy of initials, finals, and lexical tones, as well as the word
160 production intelligibility?
- 161 c) If significant improvements were observed in MMLI, could the benefits be retained after
162 the cessation of daily treatment sessions and generalize to untrained items?

163

164 **2 Methods**

165 **2.1 Participants**

166 Upon entering the treatment sessions, 30 participants with autism were randomly
167 assigned to one of two treatment groups. The MMLI group ($n = 15$, two girls) represented the
168 experimental group, and the SRT group ($n = 15$, two girls) acted as the active control group.
169 They were recruited from *Cangzhou Research Centre for Child Language Rehabilitation*.
170 Permission to conduct this study was obtained from by the local institutional review board of the
171 *Hong Kong Polytechnic University*, ensuring appropriate adherence to informed consent
172 procedures. The parents provided informed, written permission for their children to participate
173 before any assessments were administered.

174 The clinical diagnosis of ASD was established according to the DSM-5 criteria for ASD
175 (American Psychiatric Association, 2013), and further confirmed (16 out of 30 participants)
176 using the Autism Diagnostic Observation Schedule-2 (ADOS-2; Lord et al., 2012) by
177 pediatricians and child psychiatrists in child hospitals. The clinical diagnosis of ADOS-2 was not
178 conducted at the time of data collection, and the average time between the study procedures and

179 the diagnosis was 8.14 months ($SD = 5.29$ months). Since the Mandarin version of ADOS-2 has
180 not been officially validated and widely adopted in China (Sun et al. 2013; Yu et al., 2015), for
181 those cases where an ADOS report was absent (14 out of 30 participants), we confirmed
182 diagnoses using the Chinese version of the Gilliam Autism Rating Scale–Second Edition
183 (GARS-2; Gilliam, 2006) under the supervision of a licensed clinical psychologist. Other
184 inclusion criteria were: (1) the ability to imitate two piano notes; (2) the ability to sit in a chair
185 and take part in instructed activities for around 15 minutes at a time; (3) the ability to imitate
186 gross motor activities such as clapping hands, and imitate oral motor movements; (4) not having
187 the following comorbidities: cerebral palsy or tuberous sclerosis, hearing/sight impairment,
188 Down’s syndrome, uncontrolled seizures, and organic impairment of oral or laryngeal structures
189 (Chenausky et al., 2016; Wan et al., 2012). An additional 11 children with autism were found to
190 be ineligible and were excluded from the current study because they either did not meet the four
191 inclusion criteria or could not regularly attend all the required treatment sessions.

192 Furthermore, 30 TD children were recruited from one local kindergarten to investigate
193 the speech, language and cognitive levels of our participants with autism before training in
194 reference to age-matched neuro-typical ones. The TD children did not participate in the training
195 sessions. Although the TD children were not specially screened for autism symptoms via
196 standardized instruments, they met none of the diagnostic DSM-5 criteria for ASD from an
197 interview with their parents or teachers. As shown in Table 1, the participants with ASD did
198 show significant speech, language and cognitive delays compared to age-matched TD children.
199 The children with autism were divided into two subgroups in the current study: Nonverbal status
200 was defined as having complete absence of intelligible words before training; low-verbal status
201 was defined as using expressive vocabulary of no more than 50 words, based on parent reports as

202 well as the pretest measures. The term “low-verbal” was used in several previous studies (Yoder
203 & Stone, 2006; Kasari et al., 2008) to describe individuals with ASD with severe speech,
204 language and cognitive delays.

205

206 [Table 1 about here]

207

208 **2.2 Study Design**

209 **2.2.1 Language and Cognitive Evaluation**

210 A randomized control design was used in an effort to determine the effectiveness of the
211 experimental treatment (MMLI) and control for various external factors. Prior to training, all
212 participants were firstly assessed with their language ability (Ning, 2013), then nonverbal IQ
213 using the *Primary Test of Nonverbal Intelligence* (PTONI, Ehrler & McGhee, 2008), and then
214 working memory (Millman & Mattys, 2017). The two treatment groups did not differ from each
215 other in terms of chronological age, language ability, nonverbal IQ, as well as working memory
216 (Table 2).

217

218 [Table 2 about here]

219

220 *Language Ability:* The overall language ability (Chen et al., 2017; Ning, 2013) was
221 evaluated for each Mandarin-speaking participant, which consists of five subtests (including Test
222 of Mandarin Grammar, Word Definition Test, Rapid Automatized Naming, Narrative Test, and
223 Sentence Comprehension Test). These subtests evaluated both language comprehension and

224 language expression, and aimed to assess different aspects of language abilities such as
225 phonology, lexicons, grammar, and semantics. The administration time is around 30 minutes.

226 *Primary Test of Nonverbal Intelligence*: PTONI (Ehrler & McGhee, 2008) is a
227 theoretically sound, research-based method of assessing reasoning abilities in young children
228 aged from 3:0 to 9:11. The PTONI was normed on a culturally diverse demographic sample from
229 diverse language backgrounds. Testing takes approximately 15 minutes.

230 *Working Memory*: In order to evaluate the short-term phonological working memory
231 (Millman & Mattys, 2017), we administered a forward digit span task. In the task, a series of
232 numbers were played to participants auditorily and they were asked to repeat them immediately.
233 For each digit length (two to nine digits), there were two separate items. The response for each
234 item was regarded as correct and awarded 0.5 point only when the participants could correctly
235 repeat every digit in the right order. The full score for the test of forward digit span is 8.

236

237 **2.2.2 Probe Assessments**

238 For low-verbal participants, probe assessments were conducted in the Pretest (Test 1),
239 Midtest (Test 2), Immediate Posttest (Test 3), and Delayed Posttest (Test 4). After the first 12
240 treatment sessions, three nonverbal children with ASD from each treatment group (six in total)
241 demonstrated no progress at all, and still remained nonverbal. They received the second-round
242 12 treatment sessions in the same modality as the first 12 sessions, with 24 treatment sessions in
243 total (Figure 4). The change in protocol was not determined a priori. We decided to conduct
244 more training sessions for these nonverbal participants, in an effort to figure out whether the null
245 training efficacy was due to the relatively short-term training duration or due to the failure of the

246 current approach. For those six nonverbal participants, probe assessments were performed not
247 only in Tests 1–4, but also in the second-round Midtest (Test 5), and Immediate Posttest (Test 6).

248 During each probe assessment, the same picture naming task was utilized, with trained
249 (Set 1) and untrained (Set 2) picture stimuli intermixed and presented in random order. Set 1
250 consisted of 60 lexical items that were presented during probe assessments and also practiced
251 during therapy sessions. Set 2 included 12 high-frequency items, which were not practiced
252 during treatment, but presented during the probes (see Supplemental Material S1). The untrained
253 stimuli in Set 2 were used to assess the transfer of learning to untrained stimuli. The children
254 were encouraged to make more than one attempt (at most three times in total) if they did not
255 respond, or failed to produce the target word correctly, in an effort to guarantee that the
256 spontaneous productions could truly reflect production capacity in children with autism.
257 Researchers who administered the probe assessments were blind to treatment condition.
258 However, it should be noted that during the probe assessments, no correct demonstrations were
259 provided. The cueing sentence is “WHAT IS THIS?” without any information about the naming
260 of the target pictures. The spontaneous productions from each child were recorded in a sound-
261 isolated room for further analyses. As shown in Table 2, the two treatment groups did not differ
262 from each other in the production accuracy of speech sounds (initials, finals, and tones) and word
263 production intelligibility in the Pretest (i.e., at baseline assessment).

264 **2.2.3 The Development of a Smartphone/iPad App**

265 The app was developed based on the Ionic Framework
266 (<https://ionicframework.com/docs>), which uses a user interface (UI) toolkit for building mobile
267 and desktop apps via web technologies. It is an open-source front-end framework for HTML5
268 hybrid mobile application to help developers use the same source code to generate app files for

269 both Android and iOS-based platforms. The primary function of the current app is to assist
270 children’s speech learning as well as word learning. The home page presented six themes,
271 including vegetables, fruits, animals, daily necessities, snacks, and toys (Figure 2a). Under each
272 theme, a total of 10 high-frequency lexical items were chosen, resulting in 60 trained lexical
273 items in total. These 60 trained words were disyllabic nouns relevant to children’s early-acquired
274 vocabulary. These trained items contained all the 21 initials, 39 finals, four citation tones as well
275 as neutral tone, and T3 tone sandhi in Mandarin phonology. The untrained words included 12
276 high-frequency lexical items, which included early-developing ([m, n, t, t^h]), middle-developing
277 ([k, k^h, tɛ^h, ɛ]), and late-developing Mandarin consonants ([ts, ts^h, s, tʂ, tʂ^h, ʂ]), monophthongs
278 ([u, ɿ, y]), diphthongs ([ua, ou, ou]), triphthongs ([iou, iou]), finals with a nasal coda ([an, iaŋ,
279 ɔŋ, iŋ, əŋ]), as well as all the tonal categories in Mandarin. Please refer to Supplemental Material
280 S1 for details.

281

282 **[Figure 2 about here]**

283

284 Then, for each word item, one corresponding picture was presented in the center of the
285 interface as the visual cue, and two piano icons on the left and right sides (Figure 2b). The
286 natural speech sounds of 60 trained words (120 syllables) were recorded from one female
287 announcer with standard Mandarin pronunciation. To construct the piano-timbre nonspeech
288 sounds, first, a piano note (C4) of 261 Hz frequency was created, then the level pitch tier was
289 replaced with the pitch contour extracted from each syllable (120 syllables in total) using the
290 Pitch-Synchronous Overlap Add implanted in Praat (Boersma & Weenink, 2016). In this way,
291 the piano-timbre nonspeech sounds share the same pitch contours as those in natural speech

292 sounds (Figure 3). All the piano-timbre nonspeech sounds were normalized to 500 ms, and
293 equally for root-mean-square intensity level, at 70 dB SPL. During MMLI training, when tapping
294 the icons from left to right, the piano-timbre nonspeech sounds that match the pitch contours of
295 natural lexical tones for the first and second syllables would be played, one tap per syllable. As
296 for the source code, please refer to <https://github.com/introfei/VoiceTrain> and for compiling,
297 packaging, and uploading issues, please refer to <https://github.com/introfei/Blog/issues/3> to see
298 more details.

299 **[Figure 3 about here]**

300

301 **2.2.4 Treatment Protocol**

302 Therapists were female undergraduate students majoring in speech rehabilitation from the
303 *Cangzhou Research Centre for Child Language Rehabilitation*, and were trained specifically to
304 provide both MMLI and SRT training. If one therapist taught a child participant in the MMLI
305 group, she would need to teach another child participant in the SRT group. They had an average
306 of 1.73 years of experience in teaching children with autism prior to this study. Therapists firstly
307 learned from the first author on how to perform and follow the treatment protocol. They were
308 also required to make a practice of teaching other children with autism who were not
309 participating in this study until they were familiar with the whole protocol. The treatment
310 protocol of MMLI and SRT is shown in Table 3. They were conducted with intensive repetition
311 in a highly structured environment. The SRT is designed to be similar in several respects to
312 conventional forms of speech therapy, while lacking the key elements of MIT (Chenausky et al.,
313 2016). While SRT also presents verbal stimuli through the app interface and contains the same
314 steps and speech outputs as MMLI, in SRT the verbal stimuli are spoken, not intoned; and there

315 is no bimanual hand tapping on piano icons. The video demonstrations of one trial for the two
316 groups (MMLI and SRT) have been shown in Supplemental Material S2. To monitor the
317 treatment fidelity, all treatment sessions were videotaped to evaluate therapists' adherence to the
318 protocol. We reviewed five videotaped sessions (41.67 % of all treatment sessions) selected at
319 random from each child (75 treatment sessions for the MMLI group, and 75 treatment sessions
320 for the SRT control group). Each treatment session contains 30 trials (10 trials x 3 repetitions),
321 and each trial contains five steps (Table 3). On all the 2,250 MMLI trials, therapists intoned the
322 target word and tapped the piano icons, and neither of the 2,250 SRT trials were intoned or
323 tapped. Furthermore, over a total of 11,250 steps assessed in MMLI group, 91 (0.81 %) had
324 repeated steps and 44 (0.39%) had omitted steps. Over the total 11,250 steps assessed in SRT
325 group, 103 (0.92 %) had repeated steps and 38 (0.34%) had omitted steps.

326

327 [Table 3 about here]

328

329 **2.2.5 Training Procedure**

330 The therapy sessions were conducted in clinical treatment rooms at *Cangzhou Research*
331 *Centre for Child Language Rehabilitation*. The child participants in both MMLI and SRT groups
332 received short-term intensive training – 12 treatment sessions 6 times per week, over a 2-week
333 period. However, after the first round of training, three nonverbal children with ASD (out of the
334 15 participants) from each treatment group demonstrated no progress at all, and still remained
335 nonverbal. These six nonverbal participants received the second-round 12 treatment sessions,
336 with 24 treatment sessions in total (Figure 4). Each treatment session began with a warm-up
337 stage, followed by 10 lexical trials (each trial was repeated three times) of one specific theme.

338 The three repetitions of each trial were blocked together, and all the five steps of each trial were
339 repeated. The sequence of steps was the same and strictly followed the protocol regardless of the
340 response from the child (i.e., if the child did not respond or did not pronounce the word
341 correctly). Each session lasted about 50–60 minutes, including breaks, which occurred every ten
342 to fifteen minutes, based on the child’s stamina. The order of six training themes within one
343 intervention phase was randomized using a Latin Square among participants. The training order
344 in the second phase is a repetition of that in the first intervention phase for each child. While
345 receiving MMLI or SRT, the child participants were not allowed to engage in any other speech
346 therapy activities in regular school programs. As required in the consent form, parents agreed to
347 withhold other interventions while testing an unproven intervention in order to satisfy scientific
348 objectives.

349

350 **[Figure 4 about here]**

351

352 **2.3 Outcome Measures**

353 **2.3.1 Speech and Word Production Measures**

354 The recorded productions were transcribed offline by Mandarin-speaking transcribers
355 who were totally blind to the current study design in order to minimize experimental bias. Before
356 transcription, an expert majoring in phonetics picked out the best sample from each child’s
357 utterances if more than one attempt was produced for one specific target word. The criterion is to
358 choose the one with higher mean accuracy in the productions of initial, final, and lexical tone. If
359 the nonspeech-like vocalization was produced for a certain target, the nonspeech-like
360 vocalization would also be selected and sent to transcription. After selection, there were in total
361 3,634 words (7,268 syllables) produced by all the child participants from both training groups

362 across all the probe assessments. The transcribers were told about the corresponding
363 requirements in advance, but without being told about the nature and purpose of the current
364 study. The transcribers needed to transcribe all the 3,634 words (7,268 syllables) and they were
365 allowed to listen to the sound stimuli as many times as they wanted until they were confident to
366 make a transcription. The order of presentation for the transcription task was randomized. Each
367 sound stimulus was transcribed once by each transcriber.

368 First, for speech production measure, the disyllabic word was split into two syllables,
369 which were transcribed separately. The 7,268 produced syllables were randomized and further
370 transcribed using the International Phonetic Alphabet (IPA) by another five trained experts
371 majoring in linguistics (mean age = 24.50 years). Especially for the tonal coding, the exact tonal
372 categories were chosen from the following descriptions: high-level tone (T1), mid-rising tone
373 (T2 or full sandhi: *T3 [35]), dipping tone (typical T3), high-falling tone (T4), low-falling tone
374 (half-T3: *T3 [21]), and the neutral tone. Such stringent measures of phonetic transcription using
375 IPA were aimed to assess the actual correct production of phonemes (Munro & Derwing, 1995)
376 in a more fine-grained and precise manner. If none of the Mandarin phonology matched the
377 transcription, the transcribers would log “none” instead. Each coder needed around 40 hours in
378 total to complete the transcription of all the initials, finals, and tones. Three outcome measures
379 were included to evaluate the speech production accuracy: Percent Initials Correct, Percent
380 Finals Correct, and Percent Tones Correct, which were calculated with the number of correctly
381 transcribed initials, finals, and lexical tones divided by the total number of syllables,
382 respectively.

383 Second, in terms of word production measure, the 3,634 produced words were
384 randomized and presented to five native speakers of Mandarin not majoring in linguistics (mean

385 age = 21.38 years) with the E-Prime 2.0 program (Psychology Software Tools Inc., USA). The
386 transcribers were not told about the target words before transcription. They were asked to write
387 down each word they heard with two Chinese characters (each Chinese character representing a
388 morpheme in the Mandarin word) one by one in a spreadsheet to evaluate word production
389 intelligibility (i.e., how well the word is understood regardless of if all phonemes were accurately
390 produced; transcribers were allowed to take their best guess). When none of the Chinese
391 characters were suitable to be used to make a transcription, transcribers would log “none”
392 instead. Each coder needed around 15 hours to complete the transcription of words. The number
393 of correctly coded characters (morphemes) was divided by the total number of characters
394 (morphemes) to yield “Percent Morphemes Correct”.

395 To assess agreement among five raters, Kendall’s Concordance Coefficient W was
396 calculated in this study (Legendre, 2005). The interrater reliability with Kendall’s coefficients of
397 0.802 for initial transcription, 0.793 for final transcription, 0.828 for tone transcription, and 0.884
398 for morpheme coding was reached, exhibiting relatively high inter-rater reliability.

399 **2.3.2 User Experience**

400 The user experience evaluation (Chen et al., 2019) was executed after the completion of
401 all the treatment sessions, which was rated using a 5-point Likert scale based on the child’s
402 training performance (5 – the highest degree, 1 – the lowest degree). Each corresponding
403 therapist rated the user experience for her own student. Such subjective observations evaluated
404 the ways in which children with ASD approached different training methods. The user
405 experience evaluation included five aspects: enjoyment, cooperation, consistency, interest, and
406 motivation. Enjoyment refers to the degree of apparent pleasure in the learning process;
407 cooperation means the degree of collaboration in learning a trial (whether the child could follow

408 all the five steps within a trial); consistency indicates the continuity of the overall coordination
409 throughout the learning process; interest refers to the degree of apparent interest in the training
410 materials; motivation represents the degree of the apparent initiative before treatment sessions
411 (whether the child appears to want to participate in training). The operational definitions of five
412 aspects are shown in Supplemental Material S3.

413 **2.4 Statistical Analyses**

414 All the statistical analyses of outcome measures were performed in R (R Core Team,
415 2014). For the analyses of production accuracy, the generalized linear mixed-effects models
416 (GLMMs) were created using the lme4 package (Bates et al., 2014). It is feasible for GLMM in
417 R to include all the item responses transcribed from five different transcribers, and have them
418 calculated within a single statistical model. In each GLMM, *treatment group* (MMLI vs. SRT),
419 *test*, and their two-way interaction acting as fixed effects. When fitting GLMMs, *participant* and
420 *item* were included as random effects. The R code for the full model: Accuracy ~ treatment
421 group * test + (1+ test | participant: treatment group) + (1+ treatment group * test | item). By-
422 participant and by-item random intercepts and random slopes for all possible fixed factors were
423 included in the full model (Barr et al., 2013), which was compared with a simplified model that
424 excluded a specific fixed factor using the ANOVA function in lmerTest package (Kuznetsova et
425 al., 2017). Moreover, post-hoc pairwise comparisons were calculated with the lsmeans package
426 (Lenth, 2016) with Tukey adjustment.

427 For the analyses of user experience, a generalized Poisson regression model (Consul &
428 Famoye, 1992) was constructed in R using a glmer function (family = 'poisson'), with *treatment*
429 *group* (MMLI vs. SRT), *aspect* (motivation, consistency, interest, cooperation, and enjoyment),
430 and their two-way interaction acting as fixed effects. The generalized Poisson regression model

431 is often a first-choice model for counts-based datasets, and has been found useful in fitting over-
432 dispersed as well as under-dispersed count data. Given that the nonverbal and low-verbal
433 participants in this study received different amounts of treatment sessions, their results were
434 reported separately.

435

436 **3 Results**

437 **3.1 Outcomes in Nonverbal Participants**

438 There were six nonverbal participants with ASD (*G101*, *G102*, and *G103* in MMLI
439 group; *G201*, *G202*, and *G203* in SRT group), who received 24 treatment sessions across four
440 intervention phases in total (Figure 4). The probe assessment data were collected six times
441 before, during, and after therapy. Only one participant with autism from the MMLI group (*G103*)
442 began to acquire some initials, finals, lexical tones, as well as words in the trained items during
443 the second-round training, while all the other five participants remained nonverbal even after 24
444 treatment sessions. Furthermore, none of the six nonverbal participants showed any improvement
445 in the untrained items after training. For the evaluation of user experience, the subject from the
446 MMLI group (*G103*) who showed gains in speech and word learning of trained items also
447 obtained relatively higher scores of user experience, especially in the aspect of enjoyment.

448

449 **3.2 Outcomes in Low-Verbal Participants**

450 In total, there were 24 low-verbal participants with ASD ($n = 12$ in each treatment group)
451 who received 12 treatment sessions across two intervention phases (Figure 4). Figure 5 shows
452 the percentage of correct productions from two treatment groups in both trained and untrained
453 items. The x-axis represents the probe assessment sessions and the y-axis stands for the

454 percentage of correct initials, finals, lexical tones, and morphemes, respectively, from left to
455 right.

456 [Figure 5 about here]

457

458 3.2.1 Production Accuracy of Initials

459 The GLMM was performed on the production accuracy of initials in trained stimuli, and
460 the statistical results only showed a significant main effect of *test* ($\chi^2(3) = 1143.70, p < .001$).
461 However, the GLMM did not reveal the significant main effect of *treatment group* ($\chi^2(1) = 1.33,$
462 $p = .468$) nor the interaction effect of *treatment group* * *test* ($\chi^2(3) = 3.35, p = .341$). Further
463 examination on the effect of *test* implied that the low-verbal participants from both groups
464 showed significant improvement in producing initials in the trained items (all $ps < .001$), at
465 Midtest, Immediate Posttest, as well as at Delayed Posttest. Furthermore, the results showed that
466 the treatment methodology (MMLI vs. SRT) did not lead to outcome differences in the
467 production accuracy of initials in the trained stimuli across all the probe assessments (Figure 5a
468 & Table 4a).

469

470 [Table 4 about here]

471

472 Then, the GLMM on the accuracy of initials in the untrained items also merely showed a
473 main effect of *test* ($\chi^2(3) = 56.01, p < .001$), while the main effect of *treatment group* ($\chi^2(1) =$
474 $0.05, p = .818$) and the interaction of *treatment group* × *test* ($\chi^2(3) = 0.51, p = .918$) were not
475 significant. For both MMLI and SRT groups, as shown in Figure 5b, the number of correctly
476 produced initials of untrained stimuli increased significantly after the whole 12 treatment

477 sessions at Immediate Posttest ($\beta = -0.39$, $SE = 0.07$, $t = -6.03$, $p < .001$) and at follow-up
478 assessment at Delayed Posttest ($\beta = -0.34$, $SE = 0.07$, $t = -5.22$, $p < .001$). Moreover, the two
479 groups performed similarly in the production accuracy of initials in the untrained items across all
480 the probe assessments (Figure 5b & Table 4b).

481 3.2.2 Production Accuracy of Finals

482 For the trained items in Set 1, the GLMM model on the accuracy of finals showed a
483 significant main effect of *test* ($\chi^2(3) = 1192.78$, $p < .001$), while the main effect of *treatment*
484 *group* ($\chi^2(1) = 0.03$, $p = .861$) was not significant. There was a significant two-way interaction
485 of *treatment group* \times *test* ($\chi^2(3) = 17.61$, $p < .001$). Compared with Pretest, both treatment
486 groups made significant progress in the trained finals over the course of treatment (all $ps < .001$).
487 Furthermore, the two training groups differed after two intervention phases (Figure 5a & Table
488 5a), with the MMLI group showing a higher production accuracy of trained finals than the
489 matched control group at Immediate Posttest ($\beta = 0.18$, $SE = 0.04$, $t = 4.88$, $p < .001$), while no
490 group differences were found at the other probe assessments (all $ps > .05$).

491

492 [Table 5 about here]

493

494 For the untrained items in Set 2 (Figure 5b), the GLMM on production accuracy of finals
495 only reveal a main effect of *test* ($\chi^2(3) = 79.11$, $p < .001$). Neither the main effect of *treatment*
496 *group* ($\chi^2(1) = 0.18$, $p = .675$) nor the interaction of *treatment group* \times *test* ($\chi^2(3) = 3.06$, p
497 $= .382$) was significant. The post hoc analyses demonstrated that both treatment groups showed
498 significant progress in number of correctly produced finals in untrained items only at the
499 immediate posttest ($\beta = -0.59$, $SE = 0.07$, $t = -7.97$, $p < .001$). Moreover, the two groups

500 performed similarly in the production accuracy of finals in the untrained items (Figure 5b &
501 Table 5b).

502

503 3.2.3 Production Accuracy of Tones

504 The GLMM on the accuracy of trained tones showed significant main effects of
505 *treatment group* ($\chi^2(1) = 334.59, p < .001$) and *test* ($\chi^2(3) = 1493.43, p < .001$), as well as a
506 significant two-way interaction of *treatment group* \times *test* ($\chi^2(3) = 112.87, p < .001$). In
507 comparison to the performance in Pretest (Figure 5a & Table 6a), both MMLI and SRT groups
508 showed significant improvement at tone production of the trained items over the course of
509 treatment (all $ps < .001$). Nevertheless, the growth rate was quite different, with participants who
510 received MMLI training showing much higher production accuracies of trained tones compared
511 with those receiving SRT, at Midtest ($\beta = 0.22, SE = 0.04, t = 5.62, p < .001$), Immediate Posttest
512 ($\beta = 0.61, SE = 0.04, t = 16.26, p < .001$), as well as Delayed Posttest ($\beta = 0.45, SE = 0.04, t =$
513 $12.11, p < .001$).

514

515 [Table 6 about here]

516

517 Then, GLMM was performed on the accuracy of tones in untrained items, and the
518 statistical results exhibited significant main effects of *treatment group* ($\chi^2(1) = 40.85, p < .001$)
519 and *test* ($\chi^2(3) = 85.20, p < .001$), as well as a significant interaction of *treatment group* \times *test*
520 ($\chi^2(3) = 14.27, p < .01$). The low-verbal participants from the MMLI group had significant
521 progress in production of lexical tones in untrained items, at Immediate Posttest and Delayed
522 Posttest ($ps < .001$), while those from the SRT group showed no progress over the course of

523 treatment (all $ps > .05$). In terms of group difference (Figure 5b & Table 6b), the experimental
524 group of MMLI obtained much higher accuracy of tones in untrained stimuli after 12 treatment
525 sessions at Immediate Posttest ($\beta = 0.42$, $SE = 0.12$, $t = 3.41$, $p < .001$) and two weeks later at
526 Delayed Posttest ($\beta = 0.49$, $SE = 0.13$, $t = 3.90$, $p < .001$).

527 **3.2.4 Word Production Intelligibility**

528 First, for the trained items, the GLMM on word production intelligibility (Percent
529 Morphemes Correct) showed significant main effects of *treatment group* ($\chi^2(1) = 33.22$, p
530 $< .001$) and *test* ($\chi^2(3) = 1611.81$, $p < .001$), as well as a significant interaction of *treatment*
531 *group* \times *test* ($\chi^2(3) = 38.65$, $p < .001$), indicating that the two training groups showed different
532 trajectories of word learning in the trained items. As shown in Figure 5a, compared to the
533 baseline performance in Pretest, both MMLI and SRT groups showed noticeable improvements
534 in trained word production (all $ps < .001$) when tested at Midtest, Immediate Posttest, and
535 follow-up assessment two weeks later. In terms of group difference at different timepoints (Table
536 7a), the two treatment groups performed similarly on Percent Morphemes Correct in the Pretest
537 ($\beta = -0.015$, $SE = 0.038$, $t = -0.40$, $p = .688$) and Midtest ($\beta = 0.011$, $SE = 0.036$, $t = 0.29$, p
538 $= .772$), whereas the MMLI group produced higher accuracy in the trained morphemes than the
539 control SRT group in the Immediate Posttest ($\beta = 0.247$, $SE = 0.035$, $t = 7.03$, $p < .001$) as well
540 as Delayed Posttest ($\beta = 0.103$, $SE = 0.035$, $t = 2.94$, $p < .01$).

541

542 [Table 7 about here]

543

544 Second, for the untrained items, the GLMM on word production intelligibility (Percent
545 Morphemes Correct) revealed a significant main effect of *test* ($\chi^2(3) = 128.84$, $p < .001$), but the

546 main effect of *treatment group* did not reach significance ($\chi^2 (1) = 0.07, p = .789$). There was a
547 significant interaction of *treatment group* \times *test* ($\chi^2 (3) = 19.09, p < .001$). Post-hoc analysis
548 showed that compared to the performance in Pretest, the MMLI group produced more untrained
549 morphemes correctly in Midtest, Immediate Posttest, and Delayed Posttest ($ps < .05$). Over the
550 same probe assessments, however, the SRT group only produced more untrained morphemes in
551 the Immediate Posttest ($\beta = -0.54, SE = 0.10, t = -5.37, p < .001$) right after 12 treatment sessions
552 (Figure 5b & Table 7b). For the between-group difference, across all the probe assessments, the
553 MMLI group performed similarly to the control SRT group in terms of % Morphemes Correct in
554 the untrained words (all $ps > .05$).

555 **3.2.5 User experience**

556 The generalized Poisson regression model on scores of user experience for the low-verbal
557 children with ASD did not show significant main effects of *aspect* ($\chi^2 (4) = 8.22, p = .084$) and
558 *treatment group* ($\chi^2 (1) = 3.02, p = .082$). Moreover, the Poisson regression model did not reveal
559 the significant interaction of *treatment group* \times *aspect* ($\chi^2 (4) = 0.24, p = .993$). Thus, despite the
560 trend, the experimental group of MMLI ($M_{MMLI} = 3.42$) did not receive higher scores across
561 different aspects of user experience compared with the control group of SRT ($M_{SRT} = 2.65$).

562

563 **4 Discussion**

564 The current MMLI training app aimed to facilitate speech and word learning in tone-
565 language-speaking children with ASD. For the low-verbal participants, while both MMLI and
566 SRT groups were found to be effective in enhancing production skills, results suggested a faster
567 rate of improvement in the production of finals, lexical tones, and words in the trained items for
568 the experimental MMLI group. Moreover, the advantage of MMLI training transferred to the

569 untrained items in terms of lexical tone production. For the six nonverbal participants, however,
570 only one from the MMLI group responded to treatment in the trained items, while others from
571 both training groups showed no progress and remained nonverbal even after 24 treatment
572 sessions. We will discuss these findings hereunder.

573 **4.1 Potential Mechanisms Responsible for the Training Efficacy of MMLI**

574 MMLI resulted in greater improvements than SRT in most of the outcome measures for
575 the Mandarin-speaking participants, which largely corroborated the efficacy of the AMMT-based
576 training approach, with well-proven efficacy as a treatment for English-speaking children with
577 autism (Chenausky et al., 2016, 2017; Sandiford et al., 2013; Wan et al., 2011). The data
578 reported here further proved that the two key elements of MIT – intonation and hand tapping –
579 added greatly to MMLI’s effectiveness in children from another language system. When tapping
580 on the virtual piano presented through the app, the nonspeech sounds with piano timbre would be
581 generated in an effort to mimic the music-making activities. In contrast with previous studies
582 (Chenausky et al., 2016, 2017; Sandiford et al., 2013; Wan et al., 2011), our MMLI system did
583 not utilize real musical notes, but was modified to match various pitch variations of lexical tones
584 in Mandarin phonology. Hoelzley (1993) proposes that the unique timbre of musical instruments
585 may increase motivation and attention in children with autism. The MMLI group seemed to
586 enjoy the treatment more than the SRT group based on the clinical observation, but it was not
587 supported by the statistics.

588 Furthermore, music making through bimanual tapping on the tuned piano icons is a
589 multimodal activity that not only captures the attention in children with autism, but also possibly
590 primes and integrates the bilateral sensorimotor networks with shared motor, auditory, and visual
591 neural representations of the articulatory/hand movements (Bangert et al., 2006; Lahav et al.,

2007). In particular, it is speculative that the arcuate fasciculus (AF), a fiber bundle that connects the auditory perceptual regions in the temporal lobe with the motor-related regions in the frontal lobe (Catani et al., 2005), might be developed or reconstructed through intonation and music-making activities (Wan et al., 2010, 2011). The AF was thought to play an important role in auditory-motor mapping (Chenausky et al., 2017), and might be responsible for the bidirectional mapping between speech articulation and acoustics (Leclercq et al., 2010), as well as facilitating new word learning especially in the left bundle (López-Barroso et al., 2013). Another neural substrate likely to be engaged during music making is the putative mirror neuron system (MNS). It has been suggested that dysfunctional MNS underlies some of the speech and language deficits in individuals with ASD (Iacoboni & Dapretto, 2006). The elements in MMLI training, such as imitation, hand tapping, and synchronization, might activate brain regions that overlap with MNS, thus highlighting the potential benefits of such sensorimotor training to facilitate expressive language in developmental disorder such as autism (Overy & Molnar-Szakacs, 2009).

4.2 Training Efficacy for Low-Verbal Children with ASD

While MMLI holds promise for improving speech learning in Mandarin-speaking children with ASD in general, the effectiveness of MMLI was unbalanced among different components of syllables (i.e., initials, finals, and tones). As shown in Figure 5a, low-verbal participants with ASD receiving MMLI started to show superiority over those receiving SRT in their ability to correctly articulate Mandarin lexical tones in trained items as early as Midtest, and such advantage further expanded after 12 treatment sessions at Immediate Posttest and was maintained at Delayed Posttest. For speech production of Mandarin finals which use vowel(s) as the whole final or as the nucleus, the low-verbal MMLI participants only experienced comparatively greater improvement than the SRT participants after 12 treatment sessions, and

615 such advantage did not persist at the follow-up assessment. In terms of the speech production of
616 Mandarin initials that were composed of consonants, the two treatment groups performed
617 similarly in the trained items over all the probe assessments. Furthermore, for the untrained items
618 (Figure 5b), MMLI produced significantly greater gains merely in the lexical tone learning in
619 low-verbal children with ASD compared to the control therapy, SRT. Such generalization skills
620 would be greatly beneficial to children with ASD, who show difficulty in transferring learned
621 knowledge to a new context (Church et al., 2015; Happé & Frith, 2006). In a short conclusion, as
622 observed from the current data, the efficacy of MMLI was much higher in the training of lexical
623 tones, followed by vowels, and then consonants.

624 The greater improvement on lexical tone acquisition should not be surprising given that
625 relative to SRT, MMLI presented participants with additional information of pitch contours
626 embedded in piano-timbre nonspeech. Individuals with ASD have often demonstrated pitch or
627 melodic processing superiority in various musical and nonspeech stimuli (e.g. Gomot et al.,
628 2002; Heaton, 2005; O’Riordan & Passetti, 2006). Accumulating evidence pointed to a two-way
629 transferability of pitch expertise across domains of music and speech (as lexical tones) in neuro-
630 typical children and adults (Bidelman et al., 2013; Nan et al., 2018; Wong et al., 2007). By
631 targeting the clinical population, the current data provided the first empirical evidence of using
632 the relative strength of music, a ubiquitous nonspeech form, to compensate for the relative
633 weakness of speech sounds, especially lexical tone acquisition for tone-language-speaking
634 children with ASD. One recent study (Nan et al., 2018) demonstrated that the six months of
635 piano training not only enhanced the lexical tone discrimination, but also improved vowel and
636 consonant discrimination in 4- to 5-year-old Mandarin-speaking TD children, suggesting
637 strengthened common sound processing across domains underlying the benefits of musical

638 training. In this training study, however, we failed to detect the benefits of music-supported
639 MMLI training on the acquisition of initials (consonants). On the one hand, the shorter, weaker,
640 and more aperiodic consonants in speech sounds are likely to be impacted more in a co-occurring
641 nonspeech background than the stable, periodic components of tones and vowels. On the other
642 hand, since the syllable-initial consonants were acquired later than the vowels and tones in
643 Mandarin-speaking TD children (Hua & Dodd, 2000), the relatively short-term training duration
644 in this study may be another potential factor leading to the failure of transfer effects on the late-
645 acquired consonant production.

646 **4.3 Training Efficacy for Nonverbal Children with ASD**

647 The nonverbal participants in our study belong to the subgroup of the autism spectrum
648 with severe language/cognitive impairment. They were completely nonverbal, despite having
649 received extensive speech therapy (four to 16 months) prior to recruitment. In the current study,
650 except for one nonverbal child with ASD responding to the MMLI treatment, the other five
651 nonverbal participants could not correctly produce even one trained word after 24 treatment
652 sessions, and meanwhile, they received lower scores on user experience. It is unlikely that this is
653 due to the stringent measure of phonetic transcription, since these five nonverbal participants did
654 not even spontaneously produce any verbal attempts during probe assessments. In contrast, as
655 reported in Wan et al. (2011), the English-speaking nonverbal participants with autism who
656 received similar AMMT training began to vocalize some “word approximations” after 10-15
657 treatment sessions. It should be noted that the speech samples produced by the nonverbal
658 participants in Wan et al.’s (2011) study were imitations rather than spontaneous productions,
659 and should inform the degree and nature of progress. In our study, however, the spontaneous
660 speech samples were collected from participants in a picture naming task without cueing or

661 demonstration. Another possibility is that the nonverbal participants in our study were more
662 severely impaired in terms of language and cognitive capacity compared to those in Wan et al.'s
663 (2011) study. Given the extreme challenges these participants face, more treatment sessions or
664 some other training approaches should be delivered to nonverbal children with ASD. In some
665 cases, both parents and therapists have observed an increase in speechlike vocalizations during
666 vocal play in daily life, which might be an anecdotal evidence for speech development in these
667 nonverbal participants (Rvachew & Brosseau-Lapr e, 2012).

668 **4.4 Limitations and Clinical Implications**

669 This study has several limitations. First, as with many other studies of autism, a limitation
670 of the current training study is its small sample size, and replication in larger-scale randomized
671 studies will be necessary. Second, more treatment sessions or some other training approaches
672 should be applied to help nonverbal children with ASD improve their speech production skills.
673 Third, considering the substantial heterogeneity of the autism spectrum, there are various types
674 of speech disorder in ASD, such as motor speech disorder (dysarthria or apraxia), speech delay,
675 or combination of these (Chenausky et al., 2019). More in-depth investigation of speech therapy
676 in different subtypes would help determine whether MMLI is effective for all children with
677 autism or whether it only works well for certain subtypes. Understanding these mechanisms will
678 help tailor the interventions, select the most appropriate personalized treatment, and make
679 predictions about prognosis. Fourth, applying the current MMLI training to some other tonal
680 language speakers with ASD would be an important next step.

681 Taken together, the AMMT-based training program of MMLI, notwithstanding its
682 limitations, provided an effective training approach in accelerating the rate of speech sound
683 (especially lexical tone) and word production for Mandarin-speaking children with ASD. The

684 languages of the world exhibit great diversity. Some estimates suggest that around 60–70% of
685 the world’s languages are tonal (Yip, 2002), and more than half of the world’s population speak
686 a tone language (Fromkin, 1978). Thus, there is a high demand for MMLI, which could be
687 modified and applied to help some other tone-language-speaking children with autism beyond
688 Mandarin-speaking ones to better acquire the phonological category of lexical tones. With
689 respect to practical significance, the current MMLI approach is realized in the smartphone/iPad
690 app, which is easily accessible and has the potential to be utilized remotely in the home
691 environment as implemented by a parent or family member. This is important for speech therapy
692 for children with autism from counties with a shortage of speech-language pathologists. Finally,
693 the success of AMMT-based MMLI also lends support to the positive effects of music-supported
694 treatments in individuals with ASD (James et al., 2015; Reschke-Hernández, 2011; Salomon-
695 Gimmon & Elefant, 2019; Sharda et al., 2018).

696

697 **5 Conclusion**

698 This study compared the efficacy of MMLI, an AMMT-based treatment, with the control
699 therapy in eliciting spoken language for tone-language-speaking children with autism. Relative
700 to the control treatment, there was greater improvement in Mandarin-speaking children with
701 ASD after they received MMLI training, in terms of lexical tone, final, and word learning in the
702 trained items. Such enhanced training efficacy on lexical tone production persisted at two weeks
703 post-therapy, and even generalized to untrained items that were not practiced. The results hold
704 promise for the efficacy of MMLI to improve speech production in tone-language-speaking
705 children with autism. Because the nonverbal and low-verbal children with autism had a very
706 limited repertoire of speech sounds prior to treatment, the acquisition of speech sounds and

707 words through MMLI is an important gain that provides a foundation for subsequent speech and
708 language rehabilitation. More importantly, this study offers the first empirical evidence of the
709 advantages of utilizing musical elements to facilitate lexical tone acquisition in the clinical
710 population of ASD, which adds a new clinical perspective to our understanding of the close
711 relationship between music and speech.

712

713 **6 Acknowledgments**

714 This work was in part supported by the Major Program of National Social Science
715 Foundation of China (18ZDA293), the Natural Science Foundation of China overseas
716 collaboration grant (31728009), the General Research Fund of the Research Grants Council of
717 Hong Kong (15610321), and a key project from *Cangzhou Normal University* (xnjjw1907). We
718 sincerely thank all the students from *College of Qiyue Communication, Cangzhou Normal*
719 *University*, for their research assistance. We sincerely thank all the child participants from
720 *Cangzhou Research Centre for Child Language Rehabilitation* and *Tuofu Kindergarten* in
721 Cangzhou, and their parents, for their participation and cooperation.

722

723

724

725 **References**

- 726 Albert, M. L., Sparks, R. W., & Helm, N. A. (1973). Melodic intonation therapy for aphasia.
727 *Archives of Neurology*, 29(2), 130–131.
728 <https://doi.org/10.1001/archneur.1973.00490260074018>
- 729 American Psychiatric Association. (2013). *Diagnostic and statistical manual of mental*
730 *disorders: DSM-5* (5th ed.). Arlington, VA : American Psychiatric Publishing.
- 731 Bal, V. H., Katz, T., Bishop, S. L., & Krasileva, K. (2016). Understanding definitions of
732 minimally verbal across instruments: Evidence for subgroups within minimally verbal
733 children and adolescents with autism spectrum disorder. *Journal of Child Psychology and*
734 *Psychiatry*, 57(12), 1424-1433.
- 735 Baldwin, C. L. (2012). *Auditory Cognition and Human Performance: Research and*
736 *Applications*. CRC Press.
- 737 Bangert, M., Peschel, T., Schlaug, G., Rotte, M., Drescher, D., Hinrichs, H., Heinze, H.-J., &
738 Altenmüller, E. (2006). Shared networks for auditory and motor processing in
739 professional pianists: Evidence from fMRI conjunction. *NeuroImage*, 30(3), 917–926.
740 <https://doi.org/10.1016/j.neuroimage.2005.10.044>
- 741 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for
742 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*,
743 68(3), 255–278. <https://doi.org/10.1016/j.jml.2012.11.001>
- 744 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2014). Fitting linear mixed-effects models
745 using lme4. *ArXiv:1406.5823 [Stat]*. <http://arxiv.org/abs/1406.5823>
- 746 Bidelman, G. M., Hutka, S., & Moreno, S. (2013). Tone language speakers and musicians share
747 enhanced perceptual and cognitive abilities for musical pitch: Evidence for
748 bidirectionality between the domains of language and music. *PLOS ONE*, 8(4).
749 <https://doi.org/10.1371/journal.pone.0060676>
- 750 Boersma, P., & Weenink, D. (2016). *Praat: Doing phonetics by computer (Version 6.0. 14)*
751 *[Computer program]*. <http://www.praat.org/>
- 752 Buday, E. M. (1995). The effects of signed and spoken words taught with music on sign and
753 speech imitation by children with autism. *Journal of Music Therapy*, 32(3), 189–202.
754 <https://doi.org/10.1093/jmt/32.3.189>

755 Catani, M., Jones, D. K., & Ffytche, D. H. (2005). Perisylvian language networks of the human
756 brain. *Annals of Neurology*, *57*(1), 8–16. <https://doi.org/10.1002/ana.20319>

757 Chao, Y. R. (1930). A system of tone letters. *Le Maître Phonétique*, *45*, 24–27.

758 Chen, F., Peng, G., Yan, N., & Wang, L. (2017). The development of categorical perception of
759 Mandarin tones in four- to seven-year-old children. *Journal of Child Language*, *44*(6),
760 1413–1434. <https://doi.org/10.1017/S0305000916000581>

761 Chen, F., Wang, L., Peng, G., Yan, N., & Pan, X. (2019). Development and evaluation of a 3-D
762 virtual pronunciation tutor for children with autism spectrum disorders. *PLOS ONE*,
763 *14*(1), e0210858. <https://doi.org/10.1371/journal.pone.0210858>

764 Chen, F., Yan, N., Pan, X., Yang, F., Ji, Z., Wang, L., & Peng, G. (2016). *Impaired categorical*
765 *perception of Mandarin tones and its relationship to language ability in autism spectrum*
766 *disorders*. Proc. Interspeech 2016, 233-237. <https://doi.org/10.21437/Interspeech.2016->
767 1133

768 Chenausky, K., Brignell, A., Morgan, A., & Tager-Flusberg, H. (2019). Motor speech
769 impairment predicts expressive language in minimally verbal, but not low verbal,
770 individuals with autism spectrum disorder. *Autism & Developmental Language*
771 *Impairments*, *4*, 2396941519856333. <https://doi.org/10.1177/2396941519856333>

772 Chenausky, K., Norton, A., Tager-Flusberg, H., & Schlaug, G. (2016). Auditory-Motor Mapping
773 Training: Comparing the Effects of a Novel Speech Treatment to a Control Treatment for
774 Minimally Verbal Children with Autism. *PLOS ONE*, *11*(11), e0164930.
775 <https://doi.org/10.1371/journal.pone.0164930>

776 Chenausky, K. V., Norton, A. C., & Schlaug, G. (2017). Auditory-motor mapping training in a
777 more verbal child with autism. *Frontiers in Human Neuroscience*, *11*.
778 <https://doi.org/10.3389/fnhum.2017.00426>

779 Cheng, S. T. T., Lam, G. Y. H., & To, C. K. S. (2017). Pitch perception in tone language-
780 speaking adults with and without autism spectrum disorders. *I-Perception*, *8*(3),
781 2041669517711200. <https://doi.org/10.1177/2041669517711200>

782 Church, B. A., Rice, C. L., Dvongopoly, A., Lopata, C. J., Thomeer, M. L., Nelson, A., &
783 Mercado, E. (2015). Learning, plasticity, and atypical generalization in children with
784 autism. *Psychonomic Bulletin & Review*, *22*(5), 1342–1348.
785 <https://doi.org/10.3758/s13423-014-0797-9>

786 Consul, P. C., & Famoye, F. (1992). Generalized poisson regression model. *Communications in*
787 *Statistics – Theory and Methods*, 21(1), 89–109.
788 <https://doi.org/10.1080/03610929208830766>

789 Dawson, G., Meltzoff, A. N., Osterling, J., Rinaldi, J., & Brown, E. (1998). Children with autism
790 fail to orient to naturally occurring social stimuli. *Journal of Autism and Developmental*
791 *Disorders*, 28(6), 479–485. <https://doi.org/10.1023/A:1026043926488>

792 Ehrlert, D. J., & McGhee, R. L. (2008). *PTONI: Primary test of nonverbal intelligence*. Austin,
793 TX: Pro-Ed.

794 Ferguson, C. A., & Farwell, C. B. (1975). Words and sounds in early language acquisition:
795 English initial consonants in the first fifty words. *Language*, 51, 419–439.

796 Fromkin, V. (1978). *Tone: A linguistic survey*. New York: Academic Press.

797 Gilliam, J. E. (2006). *Gilliam Autism Rating Scale: GARS 2*. Austin, TX: PRO-ED.

798 Gomot, M., Giard, M.-H., Adrien, J.-L., Barthelemy, C., & Bruneau, N. (2002). Hypersensitivity
799 to acoustic change in children with autism: Electrophysiological evidence of left frontal
800 cortex dysfunctioning. *Psychophysiology*, 39(5), 577–584.
801 <https://doi.org/10.1017/S0048577202394058>

802 Hairston, M. (1990). Analyses of responses of mentally retarded autistic and mentally retarded
803 nonautistic children to art therapy and music therapy. *Journal of Music Therapy*, 27(3),
804 137–150. <https://doi.org/10.1093/jmt/27.3.137>

805 Happé, F., & Frith, U. (2006). The weak coherence account: Detail-focused cognitive style in
806 autism spectrum disorders. *Journal of Autism and Developmental Disorders*, 36(1), 5–25.

807 Heaton, P. (2005). Interval and contour processing in autism. *Journal of Autism and*
808 *Developmental Disorders*, 35(6), 787. <https://doi.org/10.1007/s10803-005-0024-7>

809 Heaton, P., Hermelin, B., & Pring, L. (1998). Autism and pitch processing: A precursor for
810 savant musical ability? *Music Perception: An Interdisciplinary Journal*, 15(3), 291–305.
811 <https://doi.org/10.2307/40285769>

812 Hoelzley, P. D. (1993). Communication potentiating sounds: Developing channels of
813 communication with autistic children through psychobiological responses to novel sound
814 stimuli. *Canadian Journal of Music Therapy*, 1, 54–76.

815 Hsieh, H. I. (1972). Lexical diffusion: Evidence from child language acquisition. *Glossa*, 6(1),
816 89–104.

817 Hua, Z., & Dodd, B. (2000). The phonological acquisition of Putonghua (Modern standard
818 Chinese). *Journal of Child Language*, 27(1), 3–42.

819 Iacoboni, M., & Dapretto, M. (2006). The mirror neuron system and the consequences of its
820 dysfunction. *Nature Reviews Neuroscience*, 7(12), 942–951.

821 James, R., Sigafoos, J., Green, V. A., Lancioni, G. E., O’Reilly, M. F., Lang, R., Davis, T.,
822 Carnett, A., Achmadi, D., Gevarter, C., & Marschik, P. B. (2015). Music therapy for
823 individuals with autism spectrum disorder: A systematic review. *Review Journal of*
824 *Autism and Developmental Disorders*, 2(1), 39–54. [https://doi.org/10.1007/s40489-014-](https://doi.org/10.1007/s40489-014-0035-4)
825 [0035-4](https://doi.org/10.1007/s40489-014-0035-4)

826 Kasari, C., Brady, N., Lord, C., & Tager-Flusberg, H. (2013). Assessing the minimally verbal
827 school-aged child with autism spectrum disorder. *Autism Research*, 6(6), 479–493.
828 <https://doi.org/10.1002/aur.1334>

829 Kasari, C., Paparella, T., Freeman, S., & Jahromi, L. B. (2008). Language outcome in autism:
830 Randomized comparison of joint attention and play interventions. *Journal of Consulting*
831 *and Clinical Psychology*, 76(1), 125–137. <https://doi.org/10.1037/0022-006X.76.1.125>

832 King, M. L., Takeguchi, K., Barry, S. E., Rehfeldt, R. A., Boyer, V. E., & Mathews, T. L.
833 (2014). Evaluation of the iPad in the acquisition of requesting skills for children with
834 autism spectrum disorder. *Research in Autism Spectrum Disorders*, 8(9), 1107–1120.
835 <https://doi.org/10.1016/j.rasd.2014.05.011>

836 Klinger, L., Dawson, G., & Renner, P. (2002). Autistic disorder. In *Child Psychopathology* (2nd
837 ed., pp. 409–454). New York: Guilford Press.

838 Koegel, R. L., Shirotova, L., & Koegel, L. K. (2009). Brief report: Using individualized
839 orienting cues to facilitate first-word acquisition in non-responders with autism. *Journal*
840 *of Autism and Developmental Disorders*, 39(11), 1587–1592.
841 <https://doi.org/10.1007/s10803-009-0765-9>

842 Kuhl, P. K., Coffey-Corina, S., Padden, D., & Dawson, G. (2005). Links between social and
843 linguistic processing of speech in preschool children with autism: Behavioral and
844 electrophysiological measures. *Developmental Science*, 8(1), F1–F12.
845 <https://doi.org/10.1111/j.1467-7687.2004.00384.x>

846 Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2017). lmerTest Package: Tests in
847 Linear Mixed Effects Models. *Journal of Statistical Software*, 82(1), 1–26.
848 <https://doi.org/10.18637/jss.v082.i13>

849 Lahav, A., Saltzman, E., & Schlaug, G. (2007). Action representation of sound: Audiomotor
850 recognition network while listening to newly acquired actions. *Journal of Neuroscience*,
851 27(2), 308–314. <https://doi.org/10.1523/JNEUROSCI.4822-06.2007>

852 Lau, J. C., To, C. K., Kwan, J. S., Kang, X., Losh, M., & Wong, P. C. (2020). Lifelong Tone
853 Language Experience does not Eliminate Deficits in Neural Encoding of Pitch in Autism
854 Spectrum Disorder. *Journal of Autism and Developmental Disorders*, 1-20.
855 <https://doi.org/10.1007/s10803-020-04796-7>

856 Leclercq, D., Duffau, H., Delmaire, C., Capelle, L., Gatignol, P., Ducros, M., Chiras, J., &
857 Lehericy, S. (2010). Comparison of diffusion tensor imaging tractography of language
858 tracts and intraoperative subcortical stimulations: Clinical article. *Journal of*
859 *Neurosurgery*, 112(3), 503–511. <https://doi.org/10.3171/2009.8.JNS09558>

860 Legendre, P. (2005). Species associations: The Kendall coefficient of concordance revisited.
861 *Journal of Agricultural, Biological, and Environmental Statistics*, 10(2), 226.
862 <https://doi.org/10.1198/108571105X46642>

863 Lenth, R. V. (2016). Least-squares means: The R package lsmeans. *Journal of Statistical*
864 *Software*, 69(1), 1–33. <https://doi.org/10.18637/jss.v069.i01>

865 Liu, Y., Li, J., Zheng, Q., Zaroff, C. M., Hall, B. J., Li, X., & Hao, Y. (2016). Knowledge,
866 attitudes, and perceptions of autism spectrum disorder in a stratified sampling of
867 preschool teachers in China. *BMC Psychiatry*, 16(1), 142.
868 <https://doi.org/10.1186/s12888-016-0845-2>

869 López-Barroso, D., Catani, M., Ripollés, P., Dell’Acqua, F., Rodríguez-Fornells, A., & Diego-
870 Balaguer, R. de. (2013). Word learning is mediated by the left arcuate fasciculus.
871 *Proceedings of the National Academy of Sciences*, 110(32), 13168–13173.
872 <https://doi.org/10.1073/pnas.1301696110>

873 Lord, C., Risi, S., DiLavore, P. S., Shulman, C., Thurm, A., & Pickles, A. (2006). Autism from 2
874 to 9 years of age. *Archives of General Psychiatry*, 63(6), 694–701.
875 <https://doi.org/10.1001/archpsyc.63.6.694>

876 Lord, C., Rutter, M., DiLavore, P. C., Risi, S., Gotham, K., & Bishop, S. (2012). *Autism*
877 *Diagnostic Observation Schedule: ADOS–2*. Los Angeles, CA: Western Psychological
878 Services.

879 Miller, S. B., & Toca, J. M. (1979). Adapted melodic intonation therapy: A case study of an
880 experimental language program for an autistic child. *The Journal of Clinical Psychiatry*,
881 *40*(4), 201–203.

882 Millman, R. E., & Mattys, S. L. (2017). Auditory verbal working memory as a predictor of
883 speech perception in modulated maskers in listeners with normal hearing. *Journal of*
884 *Speech, Language, and Hearing Research*, *60*(5), 1236–1245.
885 https://doi.org/10.1044/2017_JSLHR-S-16-0105

886 Mok, P., & Dellwo, V. (2008). Comparing native and non-native speech rhythm using acoustic
887 rhythmic measures: Cantonese, Beijing Mandarin and English. *Speech Prosody 2008*,
888 423–426.

889 Munro, M. J., & Derwing, T. M. (1995). Foreign accent, comprehensibility, and intelligibility in
890 the speech of second language learners. *Language Learning*, *45*(1), 73–97.
891 <https://doi.org/10.1111/j.1467-1770.1995.tb00963.x>

892 Nan, Y., Liu, L., Geiser, E., Shu, H., Gong, C. C., Dong, Q., Gabrieli, J. D. E., & Desimone, R.
893 (2018). Piano training enhances the neural processing of pitch and improves speech
894 perception in Mandarin-speaking children. *Proceedings of the National Academy of*
895 *Sciences*, *115*(28), E6630–E6639. <https://doi.org/10.1073/pnas.1808412115>

896 Ning, C. Y. (2013). *Test of language ability in Mandarin-speaking preschoolers*. Institute of
897 Linguistics, Tianjin Normal University: Tianjin University Press.

898 Norrelgen, F., Fernell, E., Eriksson, M., Hedvall, Å., Persson, C., Sjölin, M., Gillberg, C., &
899 Kjellmer, L. (2015). Children with autism spectrum disorders who do not develop phrase
900 speech in the preschool years. *Autism*, *19*(8), 934–943.
901 <https://doi.org/10.1177/1362361314556782>

902 Norton, A., Zipse, L., Marchina, S., & Schlaug, G. (2009). Melodic intonation therapy. *Annals of*
903 *the New York Academy of Sciences*, *1169*(1), 431–436. [https://doi.org/10.1111/j.1749-](https://doi.org/10.1111/j.1749-6632.2009.04859.x)
904 [6632.2009.04859.x](https://doi.org/10.1111/j.1749-6632.2009.04859.x)

905 O’Riordan, M., & Passetti, F. (2006). Discrimination in autism within different sensory
906 modalities. *Journal of Autism and Developmental Disorders*, 36(5), 665–675.
907 <https://doi.org/10.1007/s10803-006-0106-1>

908 Overy, K., & Molnar-Szakacs, I. (2009). Being together in time: Musical experience and the
909 mirror neuron system. *Music Perception*, 26(5), 489–504.
910 <https://doi.org/10.1525/mp.2009.26.5.489>

911 Peretz, I., Vuvan, D., Lagrois, M.-É., & Armony, J. L. (2015). Neural overlap in processing
912 music and speech. *Philosophical Transactions of the Royal Society B: Biological
913 Sciences*, 370(1664), 20140090. <https://doi.org/10.1098/rstb.2014.0090>

914 Plesa-Skwerer, D., Jordan, S. E., Brukilacchio, B. H., & Tager-Flusberg, H. (2016). Comparing
915 methods for assessing receptive language skills in minimally verbal children and
916 adolescents with ASD. *Autism*, 20, 591–604.

917 R Core Team. (2014). *R: A Language and Environment for Statistical Computing*. Vienna: R
918 Foundation for Statistical Computing. <http://www.R-project.org/>

919 Reschke-Hernández, A. E. (2011). History of music therapy treatment interventions for children
920 with autism. *Journal of Music Therapy*, 48(2), 169–207.
921 <https://doi.org/10.1093/jmt/48.2.169>

922 Rose, V., Trembath, D., Keen, D., & Paynter, J. (2016). The proportion of minimally verbal
923 children with autism spectrum disorder in a community-based early intervention
924 programme. *Journal of Intellectual Disability Research*, 60(5), 464–477.
925 <https://doi.org/10.1111/jir.12284>

926 Rvachew, S., Brosseau-Lapré, F. (2012). *Developmental Phonological Disorders: Foundations of
927 Clinical Practice*. San Diego, CA: Plural Publishing.

928 Salomon-Gimmon, M., & Elefant, C. (2019). Development of vocal communication in children
929 with autism spectrum disorder during improvisational music therapy. *Nordic Journal of
930 Music Therapy*, 28(3), 174–192. <https://doi.org/10.1080/08098131.2018.1529698>

931 Sandiford, G. A., Mainess, K. J., & Daher, N. S. (2013). A pilot study on the efficacy of melodic
932 based communication therapy for eliciting speech in nonverbal children with autism.
933 *Journal of Autism and Developmental Disorders*, 43(6), 1298–1307.
934 <https://doi.org/10.1007/s10803-012-1672-z>

- 935 Sharda, M., Tuerk, C., Chowdhury, R., Jamey, K., Foster, N., Custo-Blanch, M., Tan, M., Nadig,
936 A., & Hyde, K. (2018). Music improves social communication and auditory–motor
937 connectivity in children with autism. *Translational Psychiatry*, 8(1), 1–13.
938 <https://doi.org/10.1038/s41398-018-0287-3>
- 939 Sparks, R., Helm, N., & Albert, M. (1974). Aphasia rehabilitation resulting from melodic
940 intonation therapy. *Cortex*, 10(4), 303–316. [https://doi.org/10.1016/S0010-9452\(74\)80024-9](https://doi.org/10.1016/S0010-9452(74)80024-9)
- 942 Sun, X., Allison, C., Matthews, F. E., Sharp, S. J., Auyeung, B., Baron-Cohen, S., & Brayne, C.
943 (2013). Prevalence of autism in mainland China, Hong Kong and Taiwan: a systematic
944 review and metaanalysis. *Molecular Autism*, 4(1), 7. <https://doi.org/10.1186/2040-2392-4-7>
- 946 Sun, X., Allison, C., Wei, L., Matthews, F. E., Auyeung, B., Wu, Y. Y., Griffiths, S., Zhang, J.,
947 Baron-Cohen, S., & Brayne, C. (2019). Autism prevalence in China is comparable to
948 Western prevalence. *Molecular Autism*, 10(1), 7. <https://doi.org/10.1186/s13229-018-0246-0>
- 950 Tager-Flusberg, H., & Kasari, C. (2013). Minimally verbal school-aged children with autism
951 spectrum disorder: The neglected end of the spectrum. *Autism Research*, 6(6), 468–478.
952 <https://doi.org/10.1002/aur.1329>
- 953 Turner, L. M., Stone, W. L., Pozdol, S. L., & Coonrod, E. E. (2006). Follow-up of children with
954 autism spectrum disorders from age 2 to age 9. *Autism*, 10(3), 243–265.
955 <https://doi.org/10.1177/1362361306063296>
- 956 Wan, C. Y., Bazen, L., Baars, R., Libenson, A., Zipse, L., Zuk, J., Norton, A., & Schlaug, G.
957 (2011). Auditory-motor mapping training as an intervention to facilitate speech output in
958 non-verbal children with autism: A proof of concept study. *PLOS ONE*, 6(9), e25505.
- 959 Wan, C. Y., Demaine, K., Zipse, L., Norton, A., & Schlaug, G. (2010). From music making to
960 speaking: Engaging the mirror neuron system in autism. *Brain Research Bulletin*, 82(3),
961 161–168. <https://doi.org/10.1016/j.brainresbull.2010.04.010>
- 962 Wan, C. Y., Rüber, T., Hohmann, A., & Schlaug, G. (2010). The therapeutic effects of singing
963 in neurological disorders. *Music Perception: An Interdisciplinary Journal*, 27(4), 287–
964 295. <https://doi.org/10.1525/mp.2010.27.4.287>
- 965 Wang, W. S. Y. (1973). The Chinese language. *Scientific American*, 228(2), 50–60.

966 Wang, W. S. Y. (1978). The three scales of diachrony. In B. B. Kachru (Ed.), *Linguistics in the*
967 *Seventies: Directions and Prospects* (pp. 63–75). University of Illinois: Urbana, IL.

968 Wang, X., Wang, S., Fan, Y., Huang, D., & Zhang, Y. (2017). Speech-specific categorical
969 perception deficit in autism: An Event-Related Potential study of lexical tone processing
970 in Mandarin-speaking children. *Scientific Reports*, 7, 43254.
971 <https://doi.org/10.1038/srep43254>

972 Wong, P. C. M., Skoe, E., Russo, N. M., Dees, T., & Kraus, N. (2007). Musical experience
973 shapes human brainstem encoding of linguistic pitch patterns. *Nature Neuroscience*,
974 10(4), 420–422. <https://doi.org/10.1038/nn1872>

975 Wu, H., Lu, F., Yu, B., & Liu, Q. (2020). Phonological acquisition and development in
976 Putonghua-speaking children with autism spectrum disorders. *Clinical Linguistics &*
977 *Phonetics*, 34(9), 844-860.

978 Yip, M. (2002). *Tone*. Cambridge: Cambridge University Press.

979 Yoder P. & Stone W. L. (2006). A Randomized Comparison of the Effect of Two Prelinguistic
980 Communication Interventions on the Acquisition of Spoken Communication in
981 Preschoolers With ASD. *Journal of Speech, Language, and Hearing Research*, 49(4),
982 698–711. <https://doi.org/10.1044/1092-4388>

983 Yu, L., Fan, Y., Deng, Z., Huang, D., Wang, S., & Zhang, Y. (2015). Pitch processing in tonal-
984 language-speaking children with autism: An event-related potential study. *Journal of*
985 *Autism and Developmental Disorders*, 45(11), 3656–3667.
986 <https://doi.org/10.1007/s10803-015-2510-x>

987 Yu, L., Stronach, S., & Harrison, A. J. (2020). Public knowledge and stigma of autism spectrum
988 disorder: Comparing China with the United States. *Autism*, 24(6) 1531–1545.
989 <https://doi.org/10.1177/1362361319900839>

990

991

992

993

994 **Tables**

995

996 **Table 1**

997 *Characteristics of children with ASD and age-matched TD children*

	ASD (<i>n</i> = 30)		TD (<i>n</i> = 30)		<i>t</i>	<i>p</i>
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>		
Age (in months)	67.80	15.06	67.57	14.12	0.56	.578
Language Ability	31.67	26.84	88.07	8.81	-11.81	<.001
Nonverbal IQ	59.53	12.32	105.83	17.34	-13.78	<.001
Working Memory	3.90	4.39	12.23	3.65	-12.00	<.001
% Initials Correct	18.69	16.71	69.74	13.81	-14.23	<.001
% Finals Correct	17.60	15.59	66.31	13.74	-14.05	<.001
% Tones Correct	17.59	15.71	68.38	14.74	-14.24	<.001
% Morphemes Correct	22.68	20.99	73.88	14.73	-12.30	<.001

998

999

1000

1001 **Table 2**

1002 *Characteristics of participants with ASD in two treatment groups*

	MMLI (<i>n</i> = 15)		SRT (<i>n</i> = 15)		<i>t</i>	<i>p</i>
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>		
Age (in months)	66.13	15.32	69.47	15.14	-0.63	.542
Language Ability	33.07	25.35	30.27	29.08	0.47	.643
Nonverbal IQ	60.53	13.86	58.53	10.97	0.65	.528
Working Memory	3.93	4.38	3.87	4.55	0.05	.962
% Initials Correct	18.49	17.81	18.90	16.17	-0.28	.785
% Finals Correct	17.88	17.05	17.32	14.60	0.22	.830
% Tones Correct	17.90	17.71	17.28	14.05	0.16	.872
% Morphemes Correct	22.44	23.48	22.92	19.02	-0.15	.887

1003

1004

1005
 1006
 1007
 1008
 1009

Table 3

The warm-up stage at the beginning of each treatment session and the five-step structure of an MMLI trial vs. an SRT trial

	MMLI (Experimental Group)	SRT (Control Group)
Warming Up	Musical melodies without lyrics are played, and musical toys such as shaking maracas are introduced to facilitate their movements. Moreover, a rhythmic tapping of the foot is also used in time to music.	Playing checkboards without verbal or with minimally verbal instruction.
Steps	MMLI Trial	SRT Trial
1. Word Introduction	Therapist introduces the target word by showing a word picture (such as “tiger”) on the phone/iPad app and then intoning (singing) the word “[lɔu35 xu214]” by tapping the piano icons 1× per syllable.	Therapist introduces the target word by showing a word picture (such as “tiger”) on the phone/iPad screen and then speaking the word “[lɔu35 xu214]” without finger tapping.
2. Synchronous Production	Therapist produces target with the child. Therapist intones and taps “Let’s sing it together” and in synchrony with child “[lɔu35 xu214]”.	Therapist produces target with the child. Therapist speaks “Let’s speak it together” and in synchrony with child “[lɔu35 xu214]”.
3. Unison with Fading	Therapist and participant begin to intone and tap the target word together, but after the first syllable, the therapist stops while the child continues to intone and tap the next syllable. “[lɔu35] ____”.	Therapist and participant begin to speak the target word together, but after the first syllable, the therapist stops while the child continues to produce the next syllable. “[lɔu35] ____”.
4. Immediate Imitation	Therapist firstly intones and taps the target word alone. Afterwards, participant imitates the word, and therapist remains silent. “My turn first: [lɔu35 xu214]. Now your turn: _____”.	Therapist firstly speaks the target word alone. Afterwards, participant imitates the word, and therapist remains silent. “My turn first: [lɔu35 xu214]. Now your turn: _____”.
5. Independent Production	The child is further encouraged to independently intone and tap the target word once again. “_____”	The child is further encouraged to independently speak the target word once again. “_____”

1010
 1011

1012 **Table 4**

1013 *Means (and standard deviations) of % Initials Correct for low-verbal participants by treatment*
 1014 *group (MMLI and SRT) and probe assessment in (a) Trained Items, and (b) Untrained Items.*

	Group	Pretest		Midtest		Immediate Posttest		Delayed Posttest	
		<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
(a)	MMLI	23.41	11.83	26.93	11.35	38.48	10.45	36.80	10.98
Trained	SRT	24.03	13.27	27.92	12.90	38.03	11.51	38.34	12.63
(b)	MMLI	21.66	13.95	22.85	17.64	28.17	18.11	27.82	14.97
Untrained	SRT	21.94	14.47	22.45	18.62	28.42	15.18	26.91	12.55

1015

1016

1017 **Table 5**

1018 *Means (and standard deviations) of % Finals Correct for low-verbal participants by treatment*
 1019 *group (MMLI and SRT) and probe assessment in (a) Trained Items, and (b) Untrained Items.*

	Group	Pretest		Midtest		Immediate Posttest		Delayed Posttest	
		<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
(a)	MMLI	22.93	12.92	25.54	14.02	37.98	12.61	35.60	12.04
Trained	SRT	22.42	14.28	26.25	13.37	34.21	13.55	33.78	14.02
(b)	MMLI	19.54	14.31	20.36	14.62	27.71	17.58	23.51	15.46
Untrained	SRT	18.03	16.93	18.57	15.18	24.89	14.37	19.08	13.89

1020

1021 **Table 6**

1022 *Means (and standard deviations) of % Tones Correct for low-verbal participants by treatment*
 1023 *group (MMLI and SRT) and probe assessment in (a) Trained Items, and (b) Untrained Items.*

	Group	Pretest		Midtest		Immediate Posttest		Delayed Posttest	
		<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
(a)	MMLI	22.43	9.83	30.32	7.44	45.01	5.47	42.02	5.67
Trained	SRT	21.57	11.54	26.37	10.55	32.85	11.70	33.03	13.42
(b)	MMLI	22.53	15.06	26.03	14.08	34.78	14.95	34.03	14.90
Untrained	SRT	22.42	14.84	22.19	17.45	27.27	15.27	25.57	15.14

1024

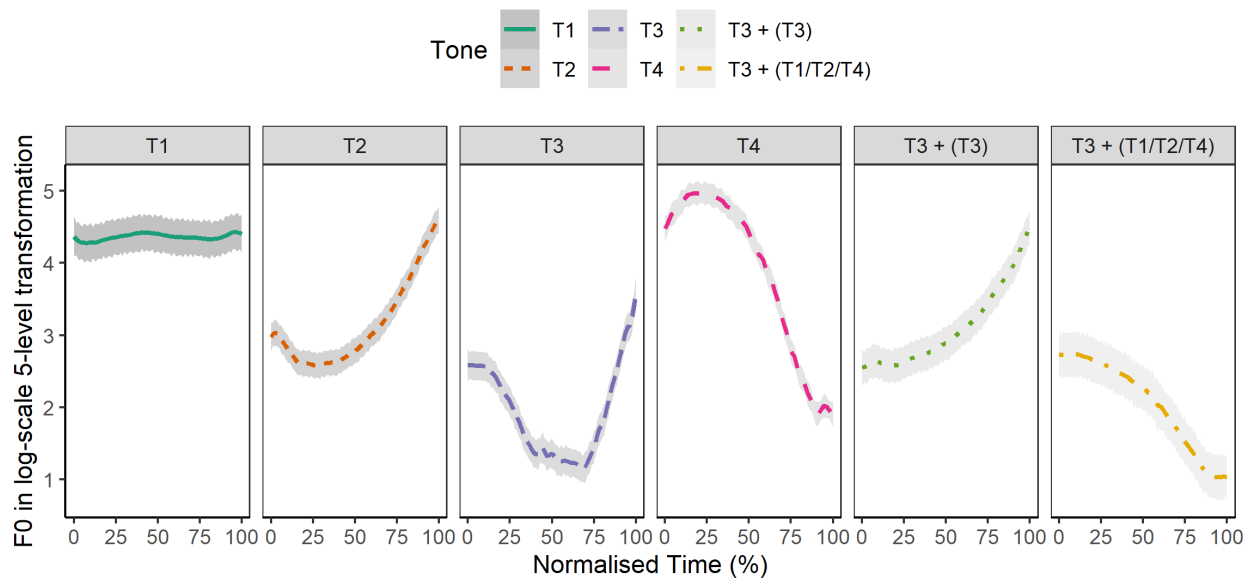
1025

1026 **Table 7**

1027 *Means (and standard deviations) of % Morphemes Correct for low-verbal participants by*
 1028 *treatment group (MMLI and SRT) and probe assessment in (a) Trained Items, and (b) Untrained*
 1029 *Items.*

	Group	Pretest		Midtest		Immediate Posttest		Delayed Posttest	
		<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
(a)	MMLI	28.71	14.03	36.17	8.22	50.87	8.00	47.72	9.18
Trained	SRT	29.05	12.64	35.83	13.20	45.32	13.30	44.35	14.82
(b)	MMLI	24.96	17.97	28.62	13.82	37.62	14.10	34.41	14.73
Untrained	SRT	27.02	15.24	26.04	19.45	34.48	17.95	29.07	17.33

1030



1031
 1032 *Figure 1.* Pitch contours of four Mandarin citation tones (T1 [55], T2 [35], T3 [214], and T4
 1033 [51]) as well as two allophonic variants of T3 due to tone sandhi (full sandhi: *T3 [35] occurring
 1034 before another T3; half-T3: *T3 [21] when before T1/T2/T4). The five-level digits in square
 1035 brackets are used to transcribe tones in Chao’s tone letters (Chao, 1930). Grey shades indicate
 1036 standard error.
 1037

(a) The home page of six themes

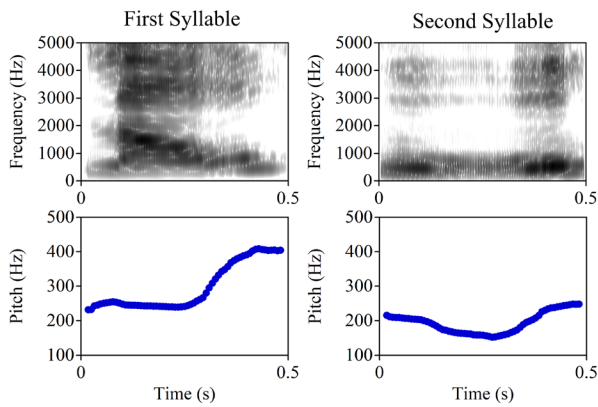


(b) The interface of one lexical item: 老虎 (tiger)

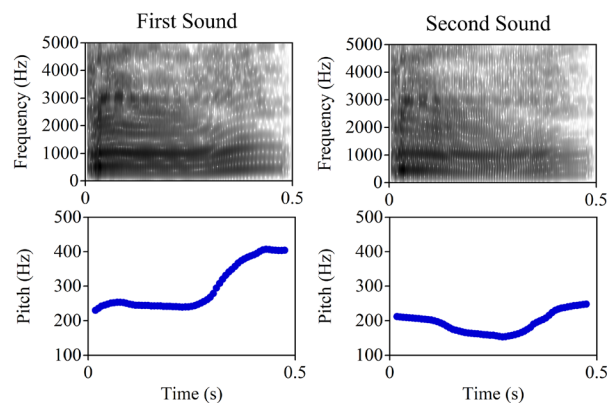


1038
 1039 *Figure 2.* The user interface of the app: (a) the home page of six themes, (b) the interface of one
 1040 lexical item.

1042 (a) Natural speech sounds: 老虎 (tiger) [lou35 xu214]



1043 (b) Corresponding piano-timbre nonspeech sounds



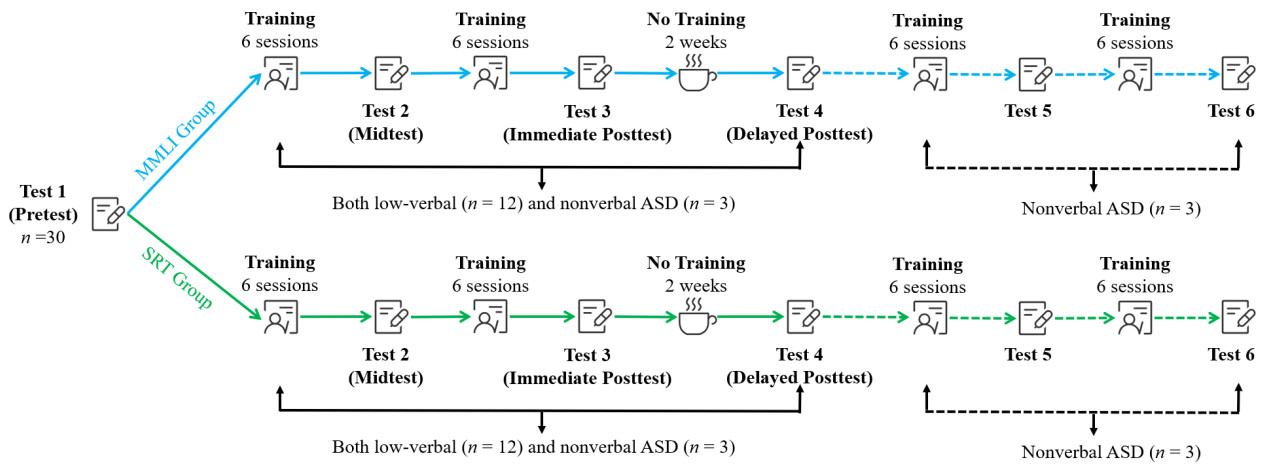
1044

1045

1046

1043 Figure 3. The spectrograms (the upper row) and pitch contours (the bottom row) of the lexical
1044 item “老虎” (tiger) [lou35 xu214] in (a) natural speech sounds, and (b) piano-timbre nonspeech
1045 sounds. The two types of sounds share exactly the same pitch contours with blue curves.

1046



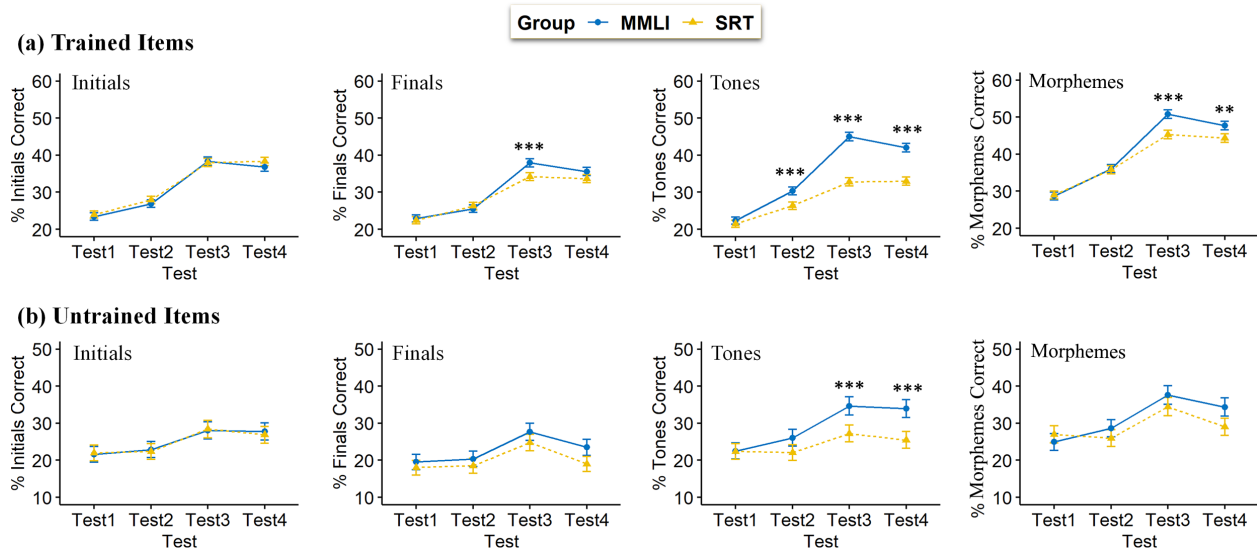
1047

1048

1049

1050

1048 Figure 4. The probe assessments and training procedure for two treatment groups (MMLI and
1049 SRT). Tests 1–6 represent the Pretest, Midtest, Immediate Posttest, and Delayed Posttest,
1050 second-round Midtest, second-round Immediate Posttest, respectively.



1051

1052 *Figure 5.* The production accuracy of initials, finals, tones, and morphemes for low-verbal
 1053 participants by treatment group (MMLI and SRT) and probe assessment in the (a) Trained Items,
 1054 and (b) Untrained Items. Tests 1–4 represent the Pretest, Midtest, Immediate Posttest, and
 1055 Delayed Posttest respectively. *** $p < .001$; ** $p < .01$ after Tukey adjustment for the comparison
 1056 of MMLI vs. SRT. Error bars: +/- 1 Confidence Interval.