

# AGREEMENT AMONG HUMAN AND AUTOMATED TRANSCRIPTIONS OF GLOBAL SONGS

**First author**

Affiliation1  
author1@ismir.edu

**Second author**

**Retain these fake authors in  
submission to preserve the formatting**

**Third author**

Affiliation3  
author3@ismir.edu

## ABSTRACT

2 Cross-cultural musical analysis requires standardized sym-  
3 bolic representation of sounds such as score notation.  
4 However, transcription into notation is usually conducted  
5 manually by ear, which is time-consuming and subjective.  
6 Our aim is to evaluate the reliability of existing methods  
7 for transcribing songs from diverse societies. We had 3 ex-  
8 perts independently transcribe a sample of 32 excerpts of  
9 traditional monophonic songs from around the world (half  
10 a cappella, half with instrumental accompaniment). 16  
11 songs also had pre-existing transcriptions created by 3 dif-  
12 ferent experts. We compared these human transcriptions  
13 against one another and against 10 automatic music tran-  
14 scription algorithms. We found that human transcriptions  
15 can be sufficiently reliable (~90% agreement,  $\kappa \sim .7$ ), but  
16 current automated methods are not (<60% agreement,  $\kappa$   
17 <.4). No automated method clearly outperformed others,  
18 in contrast to our predictions. These results suggest that  
19 improving automated methods for cross-cultural music  
20 transcription is critical for diversifying MIR.

## 1. INTRODUCTION

22 Cross-cultural analysis is essential to explore diversity and  
23 universality of music [1-2]. Such analyses require sym-  
24 bolic representations of sounds such as score notation.  
25 However, transcription into notation is usually conducted  
26 by ear, which is time-consuming and subjective [3-4].

27 Automated methods of music transcription and melody  
28 extraction might potentially solve these problems [5-7].  
29 However, automated extraction of fundamental frequency  
30 (F0) alone is not sufficient. Instead, a continuous funda-  
31 mental frequency must be segmented into discrete notes  
32 with the categorical pitches and rhythms that are distinc-  
33 tive features of almost all the world's music [8]. This chal-  
34 lenge is particularly important for variable pitch instru-  
35 ments such as the voice (the most universal instrument [8-  
36 9]). However, to our knowledge, agreement among human  
37 and automated transcription has not been objectively quan-  
38 tified using cross-cultural samples or multiple human tran-  
39 scribers.

40 The main objective of this paper is to evaluate the de-  
41 gree of agreement among human and automated transcrip-  
42 tions for a global song sample. We demonstrate that the

43 degree of agreement between human transcriptions is sub-  
44 stantially higher than the agreement between humans and  
45 machines. Our evaluation also reveals that no single algo-  
46 rithm outperforms the others, and there are no clear differ-  
47 ences between signal-processing-based methods and data-  
48 driven methods.

## 2. RELATED WORK

### 2.1 Subjectivity of manual transcription

51 Manual transcription is central to musicological research,  
52 but to our knowledge, agreement among different human  
53 transcriptions of the same songs has never been objectively  
54 measured. Even qualitative evaluation is rare. A notable  
55 exception was a 1963 symposium on transcription where  
56 four leading ethnomusicologists independently transcribed  
57 a single recording (“A Hukwe\* song with musical bow”),  
58 resulting in “four rather different transcriptions” [1, 4]. In  
59 contrast, List compared transcriptions of three songs (“Ru-  
60 manian carol”, “Yiddish lullaby”, “Thai lullaby”) by be-  
61 tween 2-9 transcribers and concluded that “transcriptions  
62 made by ear in notated form are sufficiently accurate, suf-  
63 ficiently reliable to provide a valid basis for analysis” [3].  
64 More recently, Mehr et al. [9] combined transcriptions by  
65 3 experts of 118 diverse traditional songs into a single set  
66 of “consensus” transcriptions, and had 10 experts rate their  
67 accuracy on a subjective scale from 1 (“Terrible”) to 8  
68 (“Perfect”), finding a median rating of 6 (“Very accurate”).  
69 Yet none of these studies provided an objective measure-  
70 ment of the degree of agreement between individual tran-  
71 scribers.

### 2.2 Reliability of automated transcription

73 Automatic transcription reliability has been evaluated ex-  
74 tensively for piano music and some other genres of West-  
75 ern music, but rarely for non-Western music. Ycart et al.  
76 [10] evaluated the performance of four automated tran-  
77 scription systems against perceptual ratings from 186 par-  
78 ticipants over 153 examples of piano music taken from the  
79 MAPS dataset of MIDI-aligned piano recordings [11].  
80 They found an average Fleiss’ Kappa coefficient of 0.59,  
81 or “borderline between moderate and substantial agree-  
82 ment” on participant ratings. Holzapfel and Benetos asked  
83 16 musicologists from 3 European universities to tran-  
84 scribe 8 excerpts of sousta, a traditional Greek instrumen-  
85 tal dance genre, either from scratch or starting from an au-  
86 tomatic transcription, finding “no quantitative advantage  
87 of using [automatic music transcription]” [12]. Although



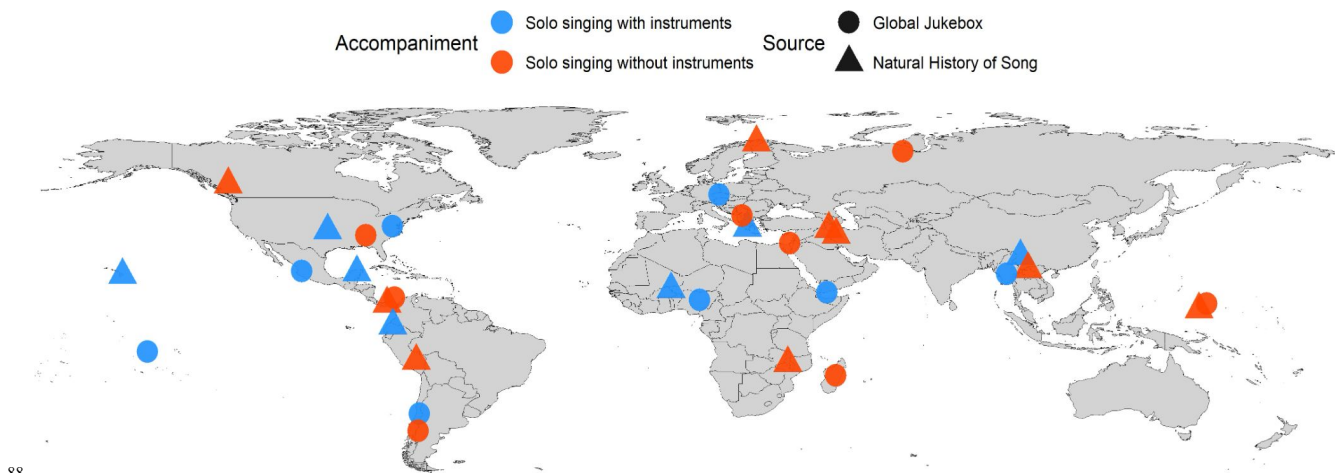


Figure 1. Map of the 32 songs transcribed and analyzed.

computer-assisted transcription studies exist [13], recent reviews by musicologists argued that computational tools for musical analysis are either useful for only low-level analysis or not widely adopted within mainstream musicology [14-15]. Overall, there is a clear need for objective measurements of agreement among automated and human transcriptions for a cross-cultural sample of songs.

### 3. METHODS AND DATASET

#### 3.1 Audio data

To examine the degree of agreement among human and automated transcriptions of diverse songs, we collected a sample of 14-second excerpts from 32 traditional songs evenly distributed across 8 geographic regions (Fig. 1). 16 songs were sampled from the publicly available 14-second excerpts of the Natural History of Song (NHS) Discography dataset [9] and manually extracted 14-second excerpts of the Global Jukebox audio files [16], respectively. We choose these datasets since they cover traditional songs from a global sample of societies. Sampling is randomly conducted using the following criteria:

- Songs are sampled equally from each of the eight regions previously used by NHS for their sampling (i.e. 4 songs per region from North America, Oceania, etc.).
- To assess capabilities of extracting vocal melodies from instrumental accompaniment, songs are sampled to consist of half solo singing without instruments and half solo singing with instruments.

One exception is that the NHS dataset contains no audio recordings of solo singing with instruments in the Middle East region, so two solo singing excerpts without instrument examples were chosen from this region instead. Sampled audio recordings contain various degrees of noise, reflecting the real-world challenges of analyzing traditional

recordings. We did not include songs with polyphonic singing since polyphonic transcription is substantially more challenging for both humans and automated methods [5], and is beyond the scope of this study.

#### 3.2 Automated methods

We selected 10 automatic music transcription/vocal melody extraction/pitch detection methods. We first choose methods listed in [10] as a baseline. However, that study focused on the systems designed for piano music, so we add methods designed for extracting pitch from human vocals. Considering the difference in the approach of the pitch estimation, our selection consists of automated methods from non-data-driven models and data-driven models. If the model employs a machine learning method (such as artificial neural networks) to learn model parameters from data in a training step, we call it data-driven, otherwise non-data-driven. Table 1 summarizes the selected automated methods. Regarding pYIN [17], we used the TONY [18] software to obtain its F0 estimation. Recently, several symbolic-level automatic transcription methods have been developed [19<sup>1</sup>-21]. However, some models were evaluated with only MIDI synthesized sounds and were not specifically designed for singing voice, so we did not select those methods.

#### 3.3 Transcription process

Twelve-tone equal temperament (12-TET) with A4 = 440 Hz is used to transcribe audio into staff notation by humans. Equal temperament is also applied to automated methods to standardize their outputs. As explained in the introduction, it is essential to obtain symbolic representations of pitch contours to analyze acoustic stimuli as melody. However, 12-TET is not completely appropriate since the pitch quantized into 12-TET does not always correspond to the actual scales/modes and perceptual tonal models even for Western singing, let alone non-Western [30-31].

While binning continuous F0 into a simplified discrete

<sup>1</sup> We have requested source code to evaluate this newly published algorithm [19] against our transcriptions, and will update our camera-ready version with this result if feasible.

Method	Target sound	Unit	Category
pYIN [17]	Monophonic vocal	Frame	Non data-driven: parameters specified manually
TONY [18]	Monophonic vocal	Note	
Melodia [22]	Vocal melody	Frame	
STF [23]	Multiple 12-tone ET	Frame	
CREPE [24]	Monophonic vocal	Frame	Data-driven: parameters optimized by training with datasets.
SPICE [25]	Monophonic vocal	Frame	
SS-nPNN [26]	Vocal melody	Frame	
AD-NNMF [27]	Multiple piano sound	Note	
OAF [28]	Multiple piano sound	Note	
madmom [29]	Multiple piano sound	Note	

**Table 1.** Summary of the selected automated methods. Unit indicates if the F0 estimation is frame-level or note-level [5] that the latter predicts onset and offset timing.

set of 12 100-cent intervals loses information about microtonal nuance, 100 cents (1 semitone) is both the most commonly used system and roughly corresponds to general levels of variability in singing intonation (imprecision and inaccuracy) [31-32], making it a reasonable choice to use to evaluate accuracy. It's also what was used by Mehr et al. [9] when creating the dataset we use, enabling us to compare our results with theirs. In summary, we decide to take advantage of the convenience and comparability of 12-TET, while acknowledging that it does not capture all musical nuances.

This study focuses on the evaluation of agreement among melodies. Hence, we discard temporal/rhythmic information and we only extract pitch from transcriptions to create a sequence of notes. However, regarding the notes representing unison melodic intervals (i.e. repeated notes), we create two transcription patterns. This is because not all selected automated methods can perform note segmentation. The change in pitch class can be used to segment two notes in the case of the other intervals, but the determining boundary between the notes of the same pitch class would require a note segmentation algorithm.

Firstly, the raw note sequences are created as a note sequence which includes the unison interval. Based on this version, we also create a note sequence which discards repeated notes and treats the notes of the unison interval as a tied single note (i.e. "CCFGGC" becomes "CFGC"). We call this version "non-unison". This treatment enables us to evaluate how much the pitch estimation itself, which is a baseline function of automatic transcription, determines

performance. In addition, 12-TET has enharmonic equivalent pitch classes, so we only use flat notes for the same sounding sharp and flat notes.

### 3.3.1 Transcription by humans

We asked three Japanese experts with professional training in Western classical music to independently transcribe the 32 recordings. One of them has professional experience of transcribing non-Western music using Western staff notation. None of them had seen the transcriptions contained in the NHS dataset. They were instructed to use MuseScore3 [33] as a tool to create transcriptions. Following Mehr et al. [9]<sup>2</sup>, we also created a consensus version of our 3 new human transcriptions. Importantly, however, while Mehr et al. only analyzed and published their consensus transcription, we include the three independent transcriptions as well as their combined consensus version to allow us to measure agreement between individual human transcribers. Our three coauthors who undertook transcription were blinded from our hypothesis testing and were asked to create transcriptions prior to discussions about coauthorship.

### 3.3.2 Transcription by automated methods

In order to standardize the output of each method, we apply post-processing steps including manual work, such as the quantization of frequency, smoothing of pitch contour, or the selection of melody contour by the Viterbi algorithm with manually specifying frequency range of melody for the case of multi-pitch estimation methods (cf. supplementary materials for details). Note that songs used in the evaluation contain solo singing with instrumental accompaniment but chosen methods are not designed to estimate the F0 of those styles of singing except for Melodia and SS-nPNN. Therefore the automated methods other than Melodia and SS-nPNN may include the pitch estimation of instrumental sounds, which is excluded from human transcriptions.

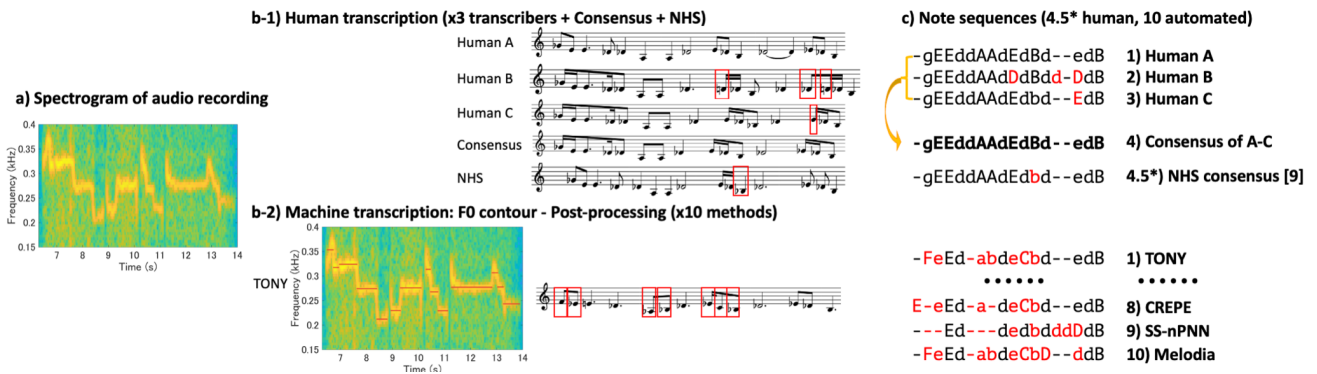
## 3.4 Sequence alignment and evaluation metrics

We use the Needleman-Wunsch algorithm [34] to align note sequences (cf. supplementary materials for further details). Agreement between two string sequences can be quantified in various ways. We mainly use Fleiss' Kappa inter-rater reliability coefficient ( $\kappa$ ), which measures how much the observed agreement exceeds chance [35]. However,  $\kappa$  does not provide other relevant information such as how many notes actually differ among note sequences or whether differences are due to disagreement about note pitch (i.e., substitution) or note segmentation (i.e., insertion/deletion). We also report such quantities by using percent identity (PID) [36-37] (cf. Fig. 2 for an example and for Supplementary Material for detailed explanation and additional analyses using Levenshtein distances).

## 3.5 Transposition

We applied transposition in the note sequence alignment process to exclude the effect of disagreement by the

<sup>2</sup> Mehr et al.'s full consensus transcriptions are published at <https://osf.io/jh7t5/>



247 **Figure 2.** Overview of the agreement evaluation using an example 8-second excerpt from NAIV-075 (Healing song,  
 248 Kwakwaka'wakw people, 00:06-00:14 from <https://osf.io/y29wp>). Red indicates disagreement with our new consensus  
 249 transcription (#4, made by combining the three individual transcriptions #1-3). For visibility, only the automated transcrip-  
 250 tion produced by TONY is shown, and octave information is omitted from the note sequences. The degree of human-human  
 251 and human-machine agreement is calculated based on the note sequences (c). For example, #5 (NHS consensus) is 95%  
 252 identical to our consensus #4 (14 out of 15 notes each), while TONY is only 48% identical (7 identical notes out of average  
 253 note length of 14.5 [38]), corresponding to Fleiss' Kappa values of .94 and .34, respectively. \*NB: NHS consensus tran-  
 254 scriptions were not available for the 16 songs from the Global Jukebox sample.

255 discrepancy of the key when calculating  $\kappa$ , PID and Le-  
 256 venshtein distance. The transposition interval was  
 257 searched from -2 semitones to +2 semitones. For human-  
 258 human transcription comparison, the transposition was  
 259 performed to maximize PID. Regarding the human-ma-  
 260 chine transcription comparison, the transposition interval  
 261 was searched to maximize the average PID of all 10 hu-  
 262 man-machine pairs for each song and each human tran-  
 263 scriber.

#### 264 4. HYPOTHESES

265 We pre-registered<sup>3</sup> the following two primary hypotheses  
 266 and 10 corresponding predictions based on pilot analysis  
 267 of 4 songs not included in our main analyses:

268 **H1: Human transcriptions are sufficiently reliable.**  
 269 This predicts a Fleiss' Kappa coefficient significantly  
 270 greater than 0 when comparing our consensus transcription  
 271 against the consensus transcriptions of Mehr et al. [9].  
 272 Note sequences including unison intervals are used.

273 **H2: TONY is the most reliable method of automated**  
 274 **singing transcription.** We predicted this because unlike  
 275 other methods TONY was designed to perform note seg-  
 276 mentation for human vocal melody, better matching hu-  
 277 man standards for transcription. This predicts that Fleiss'  
 278 Kappa comparing TONY with our consensus note se-  
 279 quences will be significantly greater than for the other 9  
 280 algorithms when evaluated against the note sequences in-  
 281 cluding unison intervals.

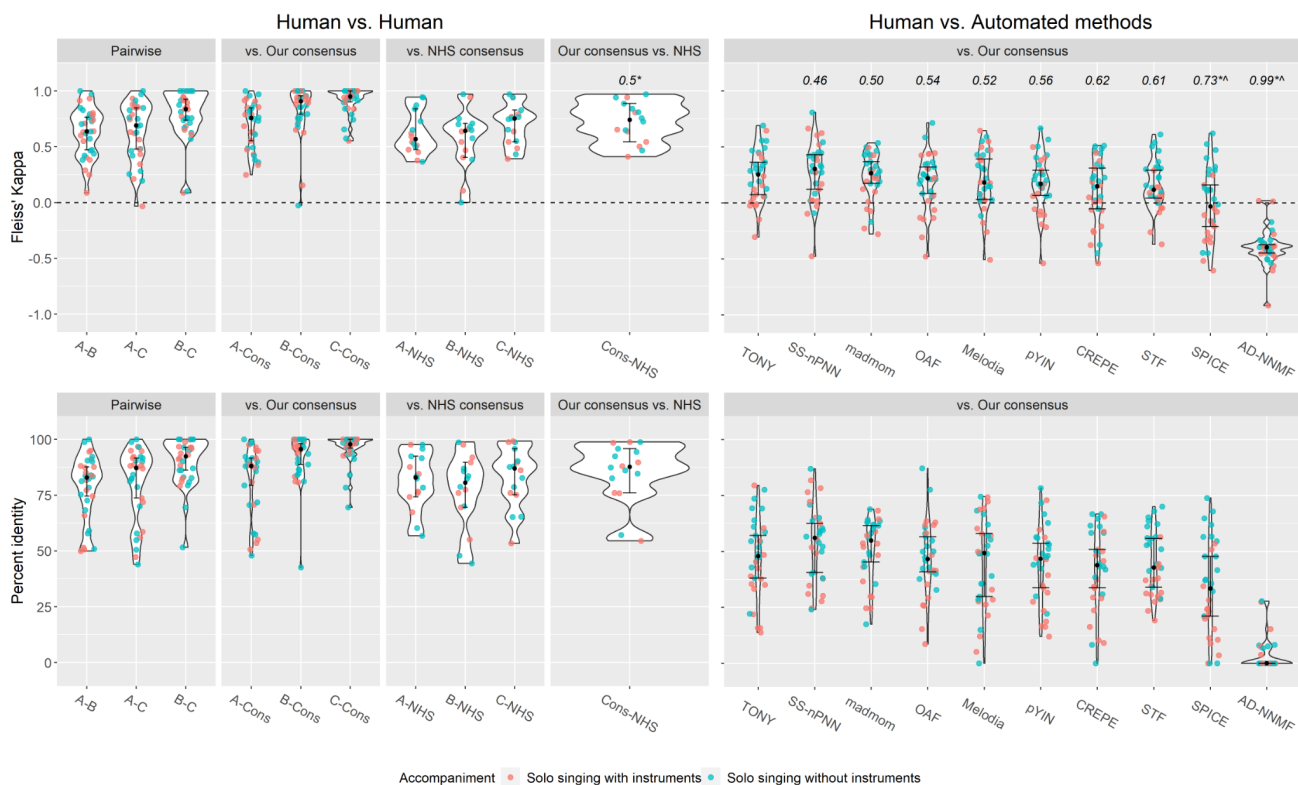
## 282 5. RESULTS

### 283 5.1 Q1. To what degree do humans' transcriptions 284 agree?

285 The left-hand side of Figure 3 shows inter-rater reliability  
 286 and percent identity results comparing human transcribers.  
 287 As predicted, there was significant agreement between our  
 288 new consensus transcriptions and the pre-existing NHS  
 289 consensus transcriptions (median  $\kappa = .74$ ,  $p < .001$ ; median  
 290 PID = 88%). When we compare the results using individ-  
 291 ual transcriptions rather than consensus transcriptions, we  
 292 see that agreement is slightly lower but still relatively high  
 293 (lowest median  $\kappa$  of .64 and PID 83% for Transcribers A  
 294 & B). The left-most two boxes show that individual vs.  
 295 consensus yields higher agreement than individual vs. in-  
 296 dividual combinations (e.g. A-Cons, A-B, A-C) for all  
 297 transcribers. This means that consensus is indeed reflect-  
 298 ing the elements of those three transcribers' transcriptions  
 299 rather than the particular pairs. These results suggest that  
 300 transcription of pitch contour could be reliable even for  
 301 non-Western music.

302 We also analyze some low agreement results. There are  
 303 25 pairwise  $\kappa$  values lower than 0.4, all of which involved  
 304 7 songs (NAIV-033, NAIV-100, NAIV-117, T5431R27,  
 305 T5482R03, NAIV-015 and NAIV-048). In particular,  
 306 NAIV-033 (Maya healing song) is a near-monotone chant,  
 307 and so the degree of agreement by chance is so large that  
 308 it negates the proportion of agreement. As the unison in-  
 309 terval dominates, the note sequences of this song are  
 310 highly homogeneous (PID > 0.9 for all 6 pairs). Other than  
 311 this song, the remaining disagreement is mainly caused by  
 312 disagreement of the pitch rather than segmentation. (cf.  
 313 supplementary material table S1). In other words, tran-  
 314 scribers generally captured the same note events, but the  
 315 assigned pitch sometimes differs by 1-2 semitones.

<sup>3</sup> <https://osf.io/bjemd> [NB: this is embargoed for anonymity; cf. Supple-  
 mentary material for an anonymized version]



316

317 **Figure 3.** Agreement among human-human and human-machine transcribed note sequences. “A”, “B”, and “C” represent  
 318 the three individual transcribers. The dashed line at  $\kappa = 0$  indicates chance levels of agreement. The numbers appearing  
 319 above the violin plots indicate effect sizes for our 10 pre-registered predictions, and \* and ^ indicate significant p-value  
 320 and posterior probability, respectively (cf. text for details). Black circles indicate medians, and bars represent 95% confi-  
 321 dence intervals of the median [38]. Results using alternative transcription methods, and full p-value and posterior proba-  
 322 bilities are available in the supplementary material (figure S3-S7 and table S2-S3).

323 Incidentally, we observed that using raw note sequences  
 324 yielded the median of  $\kappa$  very close to zero ( $\kappa = -0.019$ ) due  
 325 to cases where the tonal center differed by a semitone (and  
 326 sometimes a whole tone, cf. supplementary material figure  
 327 S1-S2 for a sample figure and the results). Therefore, our  
 328 evaluation actually focused on whether relative pitch, or  
 329 the shape of the pitch contour, matches between transcrib-  
 330 ers.

### 331 5.2 Q2. Which automated method agrees best with 332 transcription of non-Western music by humans?

333 The right-hand side of Figure 3 shows  $\kappa$  and PID obtained  
 334 by comparing the machine note sequences and our consen-  
 335 sus version’s note sequences. Contrary to our prediction,  
 336 there is no evidence for the superiority of TONY except  
 337 for AD-NNMF and SPICE. The figure also indicates gen-  
 338 erally low reliability of automated transcription methods  
 339 (median  $\kappa$  values are all below 0.4). In particular, SPICE  
 340 and AD-NNMF both had median  $\kappa$  below 0, suggesting  
 341 they performed worse than chance. Especially, AD-NNMF  
 342 failed to pick up notes correctly in many cases and indeed,  
 343 sometimes the length of note sequence of AD-NNMF is  
 344 zero (cf. supplementary material figure S8-S9). In such  
 345 cases, the proportion of agreement between human note  
 346 sequences also becomes around zero, but chance agree-  
 347 ment probability is still positive by its definition, resulting  
 348 in many negative Kappa values.

349 In addition, SPICE and CREPE had difficulty estimat-  
 350 ing F0 of the particular tracks of monophonic singing,  
 351 which is apparent from the drop in the plot of note se-  
 352 quence length (cf. supplementary material figure S8-S9).  
 353 As predicted,  $\kappa$  of automated methods designed for mono-  
 354 phonic vocal melody (i.e. TONY, pYIN, CREPE and  
 355 SPICE) show a relatively large difference dependent on in-  
 356 strumental accompaniment compared to the other methods,  
 357 but STF also suffered from instrumental sounds. (cf. Fig-  
 358 ure 3 and supplementary material figure S10).

359 See supplementary material for additional analysis de-  
 360 tails including results of measuring agreement with Le-  
 361 venshtein distances (which were generally similar to re-  
 362 sults found using PID).

## 363 6. DISCUSSION AND FUTURE WORK

364 Overall, we observed that the degree of agreement of tran-  
 365 scriptions of diverse traditional songs among human tran-  
 366 scribers was relatively high ( $\sim 90\%$  agreement,  $\kappa \sim 0.7$ ; Fig.  
 367 3), while the degree of agreement between human and au-  
 368 tomated methods was relatively low ( $< 60\%$  agreement,  $\kappa$   
 369  $< .4$ ; Fig. 3). Automated methods where less than 60% of  
 370 estimated notes agree with human judgments are unlikely  
 371 to produce satisfactory results for the kinds of tasks we  
 372 hope to use them for, such as cross-cultural comparison of  
 373 scale and interval systems [39-40]. Landis and Koch [41]  
 374 suggested that  $\kappa$  of .61-.8 be considered "substantial"

375 agreement .21-0.4 as "fair" agreement, but some have sug-  
376 gested that less than .4 is unacceptably low [42]. Our qual-  
377 itative examination of the transcriptions (e.g., Fig. 2) sup-  
378 ports the interpretation that human transcriptions of di-  
379 verse traditional songs can be sufficiently reliable, but cur-  
380 rent state-of-the-art automated methods are not.

381 Different combinations of human transcribers and  
382 songs had varying levels of agreement, but overall the  
383 agreement among three female Japanese experts and the  
384 consensus transcription by three white American male ex-  
385 perts was surprisingly high, with more differences appear-  
386 ing between individuals than between the two groups. Of  
387 course, by definition the experts had been trained in West-  
388 ern music and transcription methods - future studies should  
389 explore perceptual variability among listeners with vary-  
390 ing degrees of training in different musical systems [43].

391 Disagreement among humans appeared to primarily in-  
392 volve assignment of pitch to different pitch classes. In con-  
393 trast, disagreement in automated methods appeared to pri-  
394 marily reflect segmentation, rather than F0 estimation. Fu-  
395 ture studies might be able to clarify this point by collecting  
396 both F0 annotations and score transcriptions by humans.  
397 This might also allow us to compare our results with more  
398 conventional metrics used in research on pitch estimation  
399 algorithms such as frame-wise and note-wise F0 agree-  
400 ment and the use of true positive and false positive scores  
401 [10] (though we emphasize that our suggests that it will  
402 likely be challenging to assign a single "ground-truth" an-  
403 notation for diverse songs).

404 We were surprised that all automated methods per-  
405 formed so poorly even for the relatively simple task of  
406 transcribing only pitch sequences for monophonic songs.  
407 This suggests a strong need for automatic music transcrip-  
408 tion and other MIR tasks to expand algorithms and datasets  
409 beyond the traditional focus on Western classical and popu-  
410 lar music to be suitable for more diverse musical styles.  
411 Moving from a reliance on convenient but restricted da-  
412 taset (e.g., the MAPS dataset of MIDI-aligned piano re-  
413 cordings commonly used to evaluate automatic transcrip-  
414 tion [11]) to cross-cultural datasets like the one presented  
415 here and elsewhere [9, 16, 44] will be essential for diver-  
416 sifying MIR.

417 The formalization of a general algorithm that agrees  
418 with human pitch recognition and note segmentation is an  
419 ongoing challenge related to a central issue in MIR: the  
420 "correctness" of the algorithm depends on the degree of  
421 perceptual variability in the human ground-truth data [45].  
422 Thus, accounting for diversity and subjectivity in human  
423 transcriptions is equally critical to advance research on the  
424 automatic analysis of music. For example, while we found  
425 relatively high agreement among expert transcribers using  
426 Western 12-TET notation, we do not know whether the  
427 singers whose songs we transcribed would agree with our  
428 transcriptions, or whether transcription using a different  
429 notation system (e.g., Middle Eastern 24-note microtonal  
430 notation, 'Are'Are 7-note equiheptatonic notation [46],  
431 Killick's "global notation" [47], etc.) would give better or  
432 worse results.

433 Furthermore, here we solely focused on pitch, but a  
434 more comprehensive description of music necessitates  
435 other dimensions such as rhythm, timbre, and social con-  
436 text [48]. Other cross-cultural systems of music analysis

437 such as Cantometrics [48-49] and CantoCore [50] have  
438 been designed to capture such features. Somewhat coun-  
439 terintuitively, our current results show substantially higher  
440 agreement using Western staff notation to analyze a global  
441 song sample ( $\kappa \sim 0.7$ ) than was found using these cross-  
442 cultural song classification systems ( $\kappa \sim 0.3-0.5$  [8, 16, 50]).  
443 This suggests a need for MIR to better account for diver-  
444 sity in human ground-truth representations of all dimen-  
445 sions of music, not only pitch [51].

446 Musical diversity is a crucial challenge and opportunity  
447 for MIR. Quantifying diversity in human "ground-truth"  
448 cross-cultural data is an important first step for diversify-  
449 ing MIR. Our study demonstrates that there is still substan-  
450 tial room for improvement for automated methods of mu-  
451 sic transcription, and provides quantitative estimates of di-  
452 versity among human transcriptions to help guide future  
453 development of future MIR methods.

## 454 7. AUTHOR CONTRIBUTIONS

455 [Omitted for double-blind review]

## 456 8. ACKNOWLEDGEMENTS

457 [Omitted for double-blind review]

## 458 9. DATA/CODE AVAILABILITY

459 [Audio files, transcriptions (individual and consensus),  
460 aligned note sequences, and analysis scripts will be up-  
461 loaded to GitHub and Zenodo after blind review. Anony-  
462 mized versions can be found in the Supplementary Mate-  
463 rial]

## 464 10. REFERENCES

- 465 [1] B. Nettl, *The study of ethnomusicology: Thirty-three*  
466 *discussions, 3rd ed.*, Champaign, IL, USA: Univer-  
467 sity of Illinois Press, 2015.
- 468 [2] P. E. Savage, and S. Brown, "Toward a new compar-  
469 ative musicology," *Analytical Approaches to World*  
470 *Music*, vol. 2, no. 2, pp. 148–197, 2013.
- 471 [3] G. List, "The reliability of transcription," *Ethnomu-*  
472 *sicology*, vol. 18, no. 3, pp. 353–377, 1974.
- 473 [4] N. M. England, R. Garfias, M. Kolinski, G. List, W.  
474 Rhodes, and C. Seeger, "Symposium on transcription  
475 and analysis: A Hukwe\* song with musical bow,"  
476 *Ethnomusicology*, vol. 8, no. 3, pp. 223–233, 1964.
- 477 [5] E. Benetos, S. Dixon, Z. Duan, and S. Ewert, "Auto-  
478 matic music transcription: An overview," *IEEE Sig-*  
479 *nal Processing Magazine*, vol. 36, no. 1, pp. 20–30,  
480 2019.
- 481 [6] M. Müller, E. Gómez, and Y. Yang, Eds., "Computa-  
482 tional methods for melody and voice processing in  
483 music recordings (Dagstuhl seminar 19052)," *Dag-*  
484 *stuhl Reports*, vol. 9(1), pp. 125–177, 2019.
- 485 [7] Z. Rafii, A. Liutkus, F. Stöter, S. I. Mimilakis, D.  
486 FitzGerald, and B. Pardo, "An Overview of Lead and  
487 Accompaniment Separation in Music," *IEEE/ACM*

- 488 *Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 8, pp. 1307-1335, 2018, doi:  
489 10.1109/TASLP.2018.2825440.  
490
- 491 [8] P. E. Savage, S. Brown, E. Sakai, and T. E. Currie, "Statistical universals reveal the structures and func-  
492 tions of human music," *Proc. National Academy of*  
493 *Sciences USA*, vol. 112, no. 29, pp. 8987-8992, 2015.  
494
- 495 [9] S. A. Mehr et al., "Universality and diversity in hu-  
496 man song," *Science*, vol. 366, eaax0868, 2019, doi:  
497 10.1126/science.aax0868.
- 498 [10] A. Ycart, L. Liu, E. Benetos, and M. T. Pearce, "In-  
499 vestigating the Perceptual Validity of Evaluation  
500 Metrics for Automatic Piano Music Transcription,"  
501 *Transactions of the International Society for Music*  
502 *Information Retrieval*, vol. 3(1), pp. 68-81, 2020,  
503 doi: 10.5334/tismir.57.
- 504 [11] V. Emiya, R. Badeau, and B. David, "Multipitch es-  
505 timation of piano sounds using a new probabilistic  
506 spectral smoothness principle," *IEEE Transactions*  
507 *on Audio, Speech, and Language Processing*, vol. 18,  
508 no. 6, pp. 1643-1654, 2010.
- 509 [12] A. Holzapfel, and E. Benetos, "Automatic Music  
510 Transcription and Ethnomusicology: A User Study,"  
511 in *Proc. 20th International Society for Music Infor-*  
512 *mation Retrieval Conference*, Delft, The Netherlands,  
513 2019, pp. 678-684.
- 514 [13] E. Gómez, and J. Bonada, "Towards Computer-As-  
515 sisted Flamenco Transcription: An Experimental  
516 Comparison of Automatic Transcription Algorithms  
517 as Applied to A Cappella Singing," *Computer Music*  
518 *Journal*, vol. 37, no. 2, pp. 73-90, 2013.
- 519 [14] S. Cottrell, "Big Music Data, Musicology, and the  
520 Study of Recorded Music: Three Case Studies," *The*  
521 *Musical Quarterly*, vol. 101(2-3), pp. 216-243, doi:  
522 10.1093/musqtl/gdy013.
- 523 [15] L. Tilley, "Analytical Ethnomusicology: How We  
524 Got Out of Analysis and How to Get Back In," in  
525 *Springer Handbook of Systematic Musicology*, R.  
526 Bader, Eds. Berlin, Germany: Springer, 2018, pp.  
527 953-977.
- 528 [16] A. L. C. Wood et al., "The Global Jukebox: A public  
529 database of performing arts and culture," *PsyArXiv*  
530 *Preprint*. doi: org/10.31234/osf.io/4z97j.
- 531 [17] M. Mauch, and S. Dixon, "PYIN: A Fundamental  
532 Frequency Estimator using Probabilistic Threshold  
533 Distributions," in *Proc. IEEE International Confer-*  
534 *ence on Acoustics, Speech and Signal Processing*,  
535 Florence, Italy, 2014, pp. 659-663
- 536 [18] M. Mauch et al., "Computer-aided melody note tran-  
537 scription using the tony software: Accuracy and effi-  
538 ciency," in *Proc. 1st International Conference on*  
539 *Technologies for Music Notation and Representation*,  
540 Paris, France, 2015.
- 541 [19] R. Nishikimi, E. Nakamura, M. Goto, and K. Yoshii,  
542 "Audio-to-score singing transcription based on a  
543 CRNN-HSMM hybrid model," *APSIPA Transac-*  
544 *tions on Signal and Information Processing*, vol. 10,  
545 no. e7, pp. 1-13, 2021, doi: 10.1017/ATSIP.2021.4.
- 546 [20] R. Gunter, C. Carvalho, and P. Smaragdis, "Towards  
547 End-to-End Polyphonic Music Transcription: Trans-  
548 forming Music Audio Directly to a Score," in *IEEE*  
549 *Workshop on Applications of Signal Processing to*  
550 *Audio and Acoustics*, New York, USA, 2017, pp.  
551 151-155, doi: 10.1109/WASPAA.2017.8170013.
- 552 [21] M. A. Román, A. Pertusa, and J. Calvo-Zaragoza, "A  
553 Holistic Approach to Polyphonic Music Transcrip-  
554 tion with Neural Networks," in *Proc. 20th Interna-*  
555 *tional Society for Music Information Retrieval Con-*  
556 *ference*, Delft, The Netherlands, 2019, pp. 731-737.
- 557 [22] J. Salamon, and E. Gomez, "Melody Extraction from  
558 Polyphonic Music Signals Using Pitch Contour  
559 Characteristics," *IEEE Transactions on Audio,*  
560 *Speech, and Language Processing*, vol. 20, no. 6, pp.  
561 1759-1770, 2012, doi: 10.1109/TASL.2012.2188515.
- 562 [23] L. Su, and Y. Yang, "Combining Spectral and Tem-  
563 poral Representations for Multipitch Estimation of  
564 Polyphonic Music," *IEEE/ACM Transactions on Au-*  
565 *dio, Speech, and Language Processing*, vol. 23, no.  
566 10, pp. 1600-1612, 2015, doi:  
567 10.1109/TASLP.2015.2442411.
- 568 [24] J. W. Kim, J. Salamon, P. Li, and J. P. Bello, "Crepe:  
569 A Convolutional Representation for Pitch Estima-  
570 tion," in *Proc. 2018 IEEE International Conference*  
571 *on Acoustics, Speech and Signal Processing*, Calgary,  
572 AB, Canada, 2018, pp. 161-165, doi:  
573 10.1109/ICASSP.2018.8461329.
- 574 [25] B. Gfeller, C. Frank, D. Roblek, M. Sharifi, M. Ta-  
575 gliasacchi, and M. Velimirović, "SPICE: Self-Super-  
576 vised Pitch Estimation," *IEEE/ACM Transactions on*  
577 *Audio, Speech, and Language*, vol. 28, pp. 1118-  
578 1128, 2020, doi: 10.1109/TASLP.2020.2982285.
- 579 [26] W.-T. Lu, and L. Su, "Vocal melody extraction with  
580 semantic segmentation and audio-symbolic domain  
581 transfer learning," in *Proc. 19th International Society*  
582 *for Music Information Retrieval Conference*, Paris,  
583 France, 2018, pp. 521-528.
- 584 [27] T. Cheng, M. Mauch, E. Benetos, and S. Dixon, "An  
585 attack/decay model for piano transcription," in *Proc.*  
586 *17th International Society for Music Information Re-*  
587 *trieval Conference*, New York, NY, USA, 2016, pp.  
588 584-590.
- 589 [28] C. Hawthorne et al., "Onsets and Frames: Dual-Ob-  
590 jective Piano Transcription," in *Proc. 19th Interna-*  
591 *tional Society for Music Information Retrieval Con-*  
592 *ference*, Paris, France, 2018, pp. 50-57.
- 593 [29] S. Bock, F. Korzeniewski, J. Schluter, F. Krebs, and  
594 G. Widmer, "Madmom: A new Python Audio and  
595 Music Signal Processing Library," in *Proc. 24th*  
596 *ACM International Conference on Multimedia*, Am-  
597 sterdam, Netherlands, 2016.

- 598 [30] R. Ambrazevičius: "The Perception and Transcription of the Scale Reconsidered: Several Lithuanian Cases," *The World of Music*, vol. 47, no. 2, pp. 31-53, 2005.
- 602 [31] P. Q. Pfordresher, S. Brown, K. M. Meier, M. Belyk, and M. Liotti, "Imprecise singing is widespread," *The Journal of the Acoustical Society of America*, vol. 128(4), pp. 2182-2190, 2010.
- 606 [32] P. Larrouy-Maestri, Y. Lévêque, D. Schön, A. Giovanni, and D. Morsomme, "The Evaluation of Singing Voice Accuracy: A Comparison Between Subjective and Objective Methods," *Journal of Voice*, vol. 27, no. 2, pp. 259.e1-259.25, 2013, doi: 10.1016/j.jvoice.2012.11.003.
- 612 [33] MuseScore: <https://musescore.org/ja> (accessed Feb. 25, 2021).
- 614 [34] S. B. Needleman, and C. D. Wunsch, "A general method applicable to the search for similarities in the amino acid sequence of two proteins," *Journal of Molecular Biology*, vol. 48(3), pp. 443-453, 1970.
- 618 [35] K. L. Gwet, *Handbook of Inter-Rater Reliability. The Definitive Guide to Measuring the Extent of Agreement Among Raters, 4th edition*. Gaithersburg, MD, USA: Advanced Analytics, LLC., 2015.
- 622 [36] A. C. W. May, "Percent Sequence Identity: The Need to Be Explicit," *Structure*, vol. 12, pp. 737-738, 2004, doi: 10.1016/j.str.2004.04.001.
- 625 [37] P. E. Savage, and Q. D. Atkinson: "Automatic Tune Family Identification by Musical Sequence Alignment," in *Proc. 16th International Society for Music Information Retrieval Conference*, Málaga, Spain, 2015, pp. 162-168.
- 630 [38] M. J. Campbell, and M. J. Gardner, "Calculating confidence intervals for some non-parametric analyses," *British Medical Journal*, vol. 296, pp. 1454-1456, 1977, doi: 10.1136/bmj.296.6634.1454.
- 634 [39] J. M. McBride and T. Tlusty, "Cross-cultural data suggests musical scales evolved to maximise imperfect fifths," *arXiv Preprint*, 2020.
- 637 [40] J. Kuroyanagi *et al.*, "Automatic comparison of human music, speech, and bird song suggests uniqueness of human scales," in *Proc. 9th International Workshop on Folk Music Analysis (FMA2019)*, 2019, pp. 35-40.
- 642 [41] J. R. Landis, and G. G. Koch, "The measurement of observer agreement for categorical data," *Biometrics*, vol. 33, no. 1, pp. 159-174, 1977.
- 645 [42] J. Sim, and C. C. Wright, "The kappa statistic in reliability studies: Use, interpretation, and sample size requirements," *Physical Therapy*, vol. 85, no. 3, pp. 257-268, 2005.
- 649 [43] N. Jacoby *et al.*, "Universality and cross-cultural variation in mental representations of music revealed by global comparison of rhythm priors," *PsyArXiv Preprint*.
- 653 [44] P. E. Savage, "An overview of cross-cultural music corpus studies," in *Oxford Handbook of Music and Corpus Studies*, D. Shanahan, A. Burgoyne, and I. Quinn, Eds. Oxford University Press.
- 657 [45] A. Flexer, and T. Grill, "The Problem of Limited Inter-rater Agreement in Modelling Music Similarity," *Journal of New Music Research*, vol. 45, no. 3, pp. 239-251, 2016, doi: 10.1080/09298215.2016.1200631.
- 662 [46] H. Zemp, and V. Malkus, "Aspects of 'Are 'Are musical theory," *Ethnomusicology*, vol. 23, no. 1, pp. 5-48, 1979.
- 665 [47] A. Killick, "Global notation as a tool for cross-cultural and comparative music analysis," *Analytical Approaches to World Music*, vol. 8, no. 2, pp. 235-279, 2021.
- 669 [48] P. E. Savage, "Alan Lomax's Cantometrics Project: A comprehensive review," *Music & Science*, vol. 1, pp. 1-19, 2018.
- 672 [49] A. Lomax, and V. Grauer, "The Cantometric Coding Book," in *Folk Song Style and Culture*, A. Lomax, Ed. Washington, DC, USA: American Association for the Advancement of Science, 1968, pp. 34-74.
- 676 [50] P. E. Savage, E. Merritt, T. Rzeszutek, and S. Brown, "CantoCore: A new cross-cultural song classification scheme," *Analytical Approaches to World Music*, vol. 2, no. 1, pp. 87-137, 2012.
- 680 [51] H. Daikoku, S. Ding, U. S. Sanne, E. Benetos, A. L. Wood, S. Fujii, and P. E. Savage: "Human and automated judgements of musical similarity in a global sample," *PsyArXiv Preprint*.

## Supplementary Materials for

# AGREEMENT AMONG HUMAN AND AUTOMATED TRANSCRIPTIONS OF GLOBAL SONGS

### 1 A. PROCEDURE OF CREATING CONSENSUS NOTE SEQUENCE

2 We create a consensus transcription of each song by the following steps. Firstly, we automatically align the note sequences  
3 and then perform manual correction with rhythmic information. Secondly, disagreements of each note among note se-  
4 quences are resolved by majority rule, following [9]. If there is a note that is different in all three note sequences, we ask  
5 our collaborators via email to choose which note would fit the consensus notes selected by the majority rule. If disagreement  
6 still remains, we choose the median of the pitch from the disagreeing notes. However, subjective decisions were made  
7 when disagreement involved alignment gaps. We then asked transcribers by email to confirm the soundness of the resultant  
8 consensus transcriptions, and further updated some transcriptions based on this feedback.

### 9 B. DETAILS OF POST-PROCESSING PROCEDURE

10 In order to standardize the output of each method, we applied the following processes.

- 11
- 12 1. For methods that do not quantize F0 to twelve-tone equal temperament, the estimated F0 is rounded to the nearest  
13 frequency of the twelve-tone equal temperament.
- 14 2. For methods that do not estimate note duration or note tracking, a median filter with the length of 0.25 seconds is  
15 applied to smooth the pitch contour. Furthermore, the sequences of F0 shorter than 0.15 seconds are ignored from  
16 transcription. 0.15 is determined to make the length of note sequence similar to the humans' sequences. These param-  
17 eters were tuned to minimize the possibility that the automated methods would produce long sequences made up of  
18 unrealistically short notes as a by-product of the instability of pitch targets in human singing. If the unit of the discrete  
19 time interval of generated time-frequency representation is less than 0.01 second, decimation is applied to make the  
20 interval close to 0.01 second to smooth the pitch contour.
- 21 3. For methods predicting multiple pitches in a single timeframe, we apply the following steps to obtain the stream of  
22 single pitch prediction. Firstly, we observe that these methods tend to predict an overtone as a separate note, so the  
23 frequency range of the melody is manually specified, and the F0 prediction out of this range was removed. After that,  
24 the Viterbi algorithm is applied to the remaining multi-pitch F0 prediction results to obtain the dominant time-fre-  
25 quency energy sequence as a melody [50].
- 26 4. Regarding CREPE, F0s having a confidence score larger than or equal to 0.8 are picked up. Note that there is no  
27 guideline of what value to be used for a threshold. If we used a lower threshold, the final note sequence would become  
28 longer due to including more pitches, and that would result in lower PID and Kappa than the currently presented  
29 results.
- 30 5. We use a song excerpt as the input of automated methods to obtain the pitch estimation of a specified 14-second  
31 segment. However, pitch estimation process would depend on the information available on the broader time range of  
32 audio data to estimate the F0 of local time-frame, so feeding an entire song as input and extracting the target segment  
33 from its output will produce different pitch estimation results. In this study, we only have the excerpt of songs regard-  
34 ing the NHS, so we decided to consolidate the input by an audio excerpt.

35  
36 Incidentally, OAF estimates onset and offset of note, but it is fairly precise, so the above post-processing is applied to  
37 make a more meaningful comparison with the other methods.

### 38 C. SEQUENCE ALIGNMENT METHOD

39 We perform pairwise alignment to create the alignment of note sequences by the Needleman-Wunsch algorithm using 0.0  
40 for gap opening penalty, -1.0 for gap extension penalty and -1.0 for mismatch (substitution) penalty. This is a linear gap  
41 setting, and we choose this setting that makes the alignment score equivalent to Levenshtein distance whose operations (i.e.  
42 insertion, deletion, substitution) are all equally weighted. We use octave information for the evaluation, so the element of  
43 the sequence consists of two characters: pitch class and octave level (e.g. "A4"). When multiple sequence alignment is  
44 necessary for creating the baseline of consensus note sequences, we use the center star method to solve alignment heuris-  
45 tically since the computation of the global optimal multiple sequence alignment is not feasible due to its computational  
46 complexity [51-52]. The center sequence is determined by the sum-of-pairs scoring [51-52], and each score is calculated  
47 by the Needleman-Wunsch algorithm as described above.

### 48 D. METRICS OF AGREEMENT AMONG SEQUENCES

49 Regarding the computation of Fleiss' Kappa, we regard note transcriptions as a transcriber's categorization of the F0 of a  
50 given note. We do not apply partial agreement in this study. The length of the sequence corresponds to the number of  
51 subjects, and the number of sequences corresponds to the number of raters. When calculating the inter-rater reliability  
52 coefficients, we also treat gaps inserted during the alignment as coding rather than absence. Gap insertion indicates that  
53 some transcribers interpret the sound as a pitch, but the others do not, which we treat as a coding disagreement.

54 On the other hand, the practice of reporting the raw percent agreement score along with inter-rater reliability coefficients  
55 is also discussed due to its simplicity [35, 53]. Percent identity (PID) measures the proportion of concordance elements of  
56 two sequences which is conceptually equivalent to percent agreement, and we use this metric to evaluate how much two  
57 note sequences are identical. In the case of multiple sequences appearing in group-wise agreement evaluation, we average  
58 the PID by all combinations of pairs in the sequences. PID has originally been used in the computational analysis of protein  
59 and DNA sequences to express the similarity between two sequences [36, 52, 54], as well as the comparative study of  
60 traditional music melodies [37]. There are several variations in PID [36], and we use the following version.

$$\begin{aligned}
 \text{PID} &= 100 \left\{ \frac{N_{\text{ID}}}{0.5(L_1 + L_2)} \right\} & (1) \\
 N_{\text{ID}} &:= \text{Number of identical notes} \\
 L &:= \text{Length of sequence}
 \end{aligned}$$

65 Although Kappa coefficients and PID can provide the reliability of agreement and the proportion of equality of note  
66 sequences respectively, these quantities do not tell how many notes actually differ between note sequences. Therefore, we  
67 also use Levenshtein distance to quantify such difference by the number of insert/delete/substitution operations. The pen-  
68 alty of each operation is equally weighted by 1. The score is also averaged in groupwise evaluation cases as well as PID.

## 69 E. STATISTICAL ASSUMPTIONS OF THE TESTS

70 Inter-rater reliability coefficients, PID, and Levenshtein distance all quantify the degree of concordance among sequences.  
71 The underlying distribution of inter-rater reliability coefficients is considered to depend on the raters (i.e. transcribers) and  
72 subjects (audio recording) [35]. Furthermore, our agreement metrics are collected from various combinations of transcrib-  
73 ers and audio samples, and the domain of Kappa is finitely bounded, so the resultant distributions of agreement metrics  
74 would not necessarily fit normal or location-scale family distributions.

75 Based on the above assumption, we consider the appropriate testing methods to handle the metrics to be nonparametric  
76 methods. We choose the sign test for one-sample test case and the two-sample Anderson-Darling test [55] and two-sample  
77 Bayesian nonparametric testing using Pólya trees [56] for two-sample test scenarios. Regarding the two-sample test, we  
78 assess the probability of type I error by the two-sample Anderson-Darling test. Besides, to complement the lack of infor-  
79 mation about how much we can be confident in accepting alternative hypotheses, we also employ Bayesian hypothesis  
80 testing. Although these two tests are different procedures, both are proved to be asymptotically consistent under the null  
81 hypothesis ( $F(x) = G(x)$ ) and the alternative hypothesis ( $F(x) \neq G(x)$ ) [55-56]. Please refer to the next section for the  
82 detailed setting of Bayesian nonparametric testing.

83 Regarding the effect size to be used for our nonparametric tests, we choose the departure from the expected proportion  
84 under the null hypothesis proposed by Cohen [57] for the one-sample test and the probability-based effect size measure  $A$   
85 which is known as the probability of one group's superiority over another for two-sample tests [58]. The departure effect  
86 size (or Cohen's  $g$ ) in our study can be interpreted as follows. The sign test uses the number of samples whose value is  
87 larger than the expected median under the null hypothesis as test statistics. If the null hypothesis of the sign test is true,  
88 then the proportion of data (i.e.  $\kappa$  in our case) above the expected median (i.e. 0 in our case) should be around 50% of all  
89 samples. However, if the actual median is larger than 0, then the proportion of samples above the expected median would  
90 be larger than 0.5. We calculate the proportion of samples larger than 0 and show the difference between that proportion  
91 and the expected proportion under the null hypothesis (i.e. 0.5). Note that in this case, the range of the effect size is from 0  
92 to 0.5 and Cohen [59] suggests interpreting the value larger than 0.25 as the existence of a "Large" effect.

93 The probability-based effect size uses empirical distributions of data to quantify how much data in a group takes a larger  
94 value than another group, and it is robust to violations of the parametric assumptions. We use this effect size to interpret  
95 how much TONY's  $\kappa$  is large compared to the others. Note that  $A$  can be converted to a common standardized mean dif-  
96 ference such as Cohen's  $d$  if the normality assumption of data holds [58].

97 In summary, we put non-normality assumptions for the distributions of  $\kappa$ . Thus, we chose testing methods including  
98 Bayesian tests and effect size from nonparametric techniques. We performed the one-tailed one-sample sign test assuming  
99 the median of Fleiss' Kappa to be 0 as a null hypothesis for the hypothesis testing of human-human agreement evaluation.  
100 Regarding the hypothesis testing of examining the automated method producing transcriptions that best agree with humans'  
101 transcriptions, the null hypothesis to be tested is  $F_{\text{TONY}}(\kappa) = G_{\text{OTHER}}(\kappa)$ , which is the 9 two-sample tests of comparing the  
102 empirical distribution of  $\kappa$  by TONY and the others. The superiority of TONY can be quantified by whether the probability-  
103 based effect size measure  $A$  exceeds 0.5 or not.

104 **F. SETTING OF BAYESIAN NONPARAMETRIC TESTING**

105 We set  $c = 1$  and the normal distribution as the centering distribution as the parameters of the Pólya trees (see [56] for the  
 106 definition of parameterization of this test). However, we use the mean and standard deviation to create partitions of samples  
 107 instead of the median and quantiles used in the original study. We set the equal probability for the null hypothesis and the  
 108 alternative hypothesis ( $= 0.5$ ) as the prior distribution of our Bayesian hypothesis testing, so the posterior odds are equal to  
 109 Bayes factor.

110 **G. CONTROL OF SIMULTANEOUS INFERENCE**

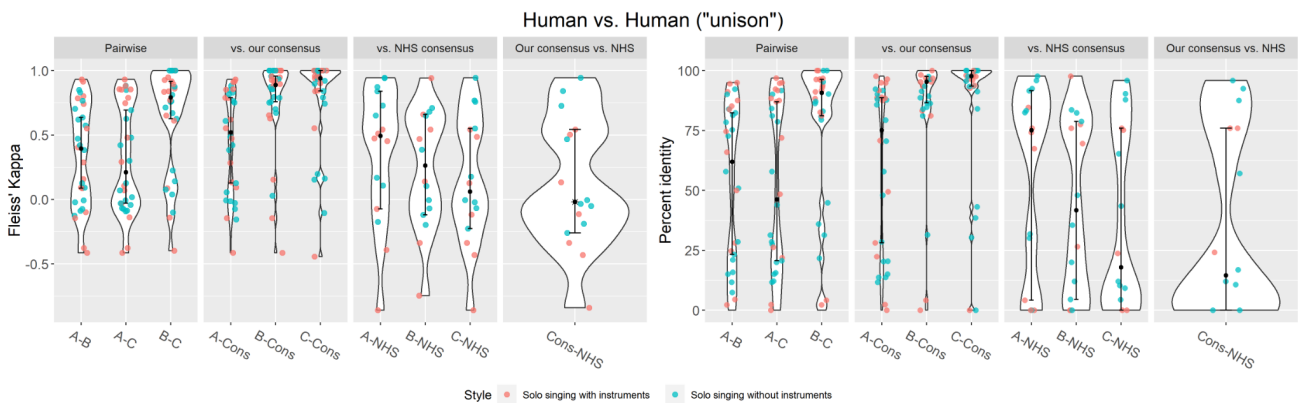
111 There are 10 null hypothesis significance tests in our analysis: one-sample Sign test  $\times 1$  + two-sample Anderson-Darling  
 112 test  $\times 9$  (machine pairs). Since our discussion on the reliability of transcription is interrelated to these test results, we use  
 113 the False Discovery Rate method to control the p-value threshold for all hypothesis tests regarding these as multiple testing  
 114 and simultaneous inference. In particular, we will use the Benjamini–Hochberg step-up procedure [60] at level  $\alpha = 0.05$  as  
 115 the threshold to determine the rejection of 10 null hypothesis testing. Incidentally, we will interpret that the Bayesian test  
 116 at least substantially supports the alternative hypothesis if the posterior probability exceeds 0.8 which corresponds to the  
 117 Bayes factor = 4 in our setting (i.e. the prior distribution being equally weighed to the null and alternative hypothesis).

118 **H. SEMITONE DISCREPANCY**

BeeeeBdBeeBBeeBBBeeeeBdBeeBBeeBBeeeeBdBBe-  
 bDDDDbCbbDDbbDDbbDDDDbCbbDDbbDDbbDDDDbCbbGC

119 **Figure S1.** Example of semitone discrepancy (NAIV-100). Octave information is omitted for visibility.  
 120

121 **I. RESULTS OF AGREEMENT USING RAW NOTE SEQUENCES**



122 **Figure S2.** Agreement of human transcriptions not applying transposition.  
 123

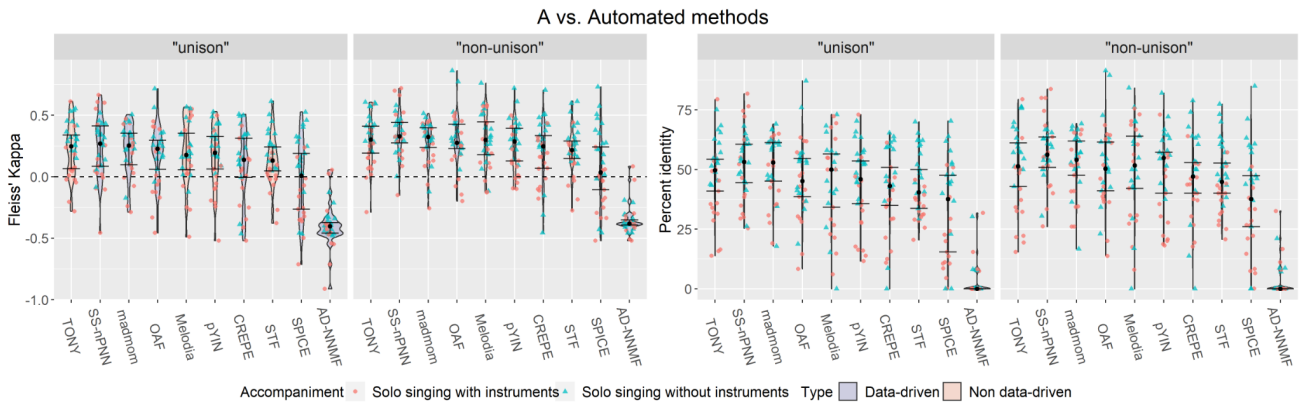
124 **J. SUMMARY OF DISAGREEMENT FACTORS OF LOW DISAGREEMENT SONGS**

125

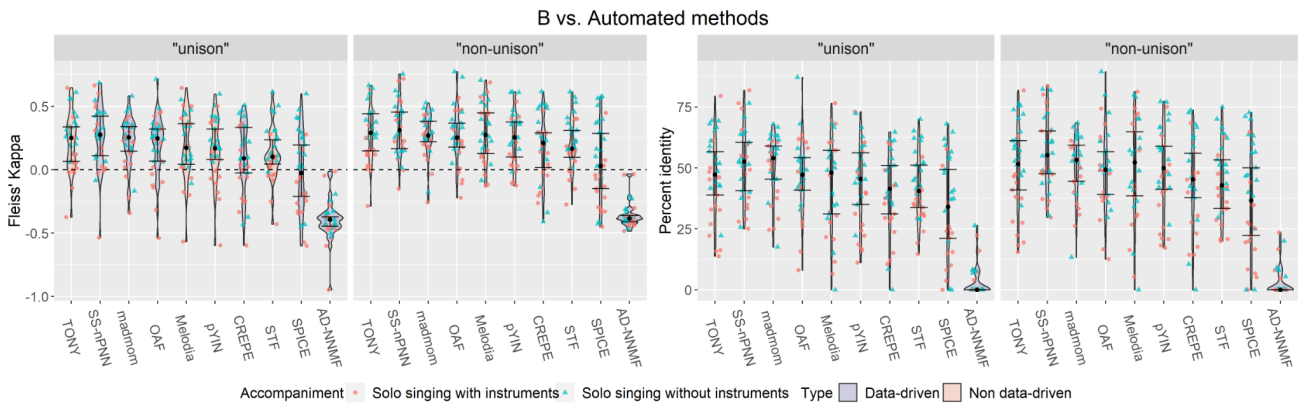
Song	Segmentation	Pitch	Both
NAIV-015	1	0	1
NAIV-048	0	0	1
NAIV-100	3	3	0
NAIV-117	0	4	0
T5431R27	0	0	3
T5482R03	0	3	0

126 **Table S1.** Qualitative classification of major disagreement factors of 19 pairs. The number indicates the count by seg-  
 127 mentation disagreement, pitch disagreement or both factors.

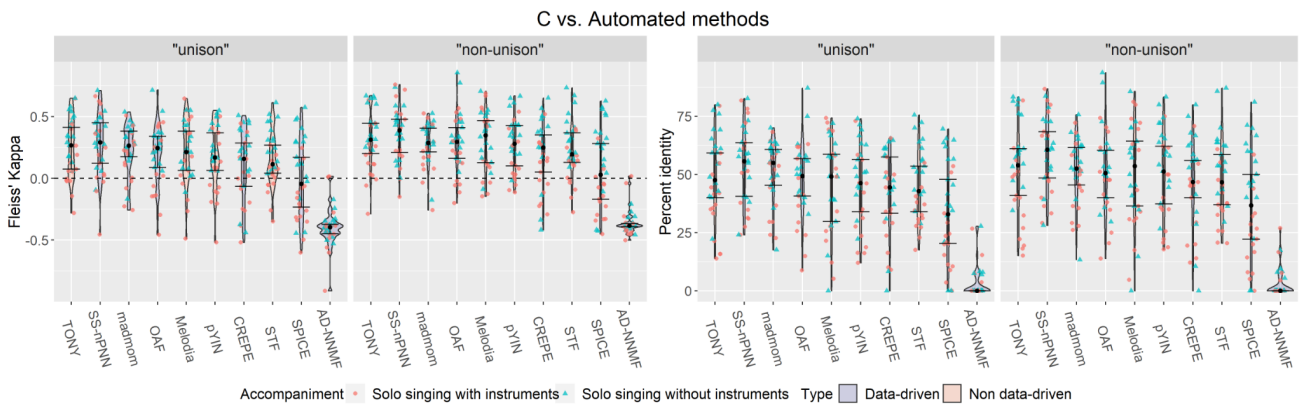
128 **K. AGREEMENT BETWEEN AUTOMATED METHODS AND INDIVIDUAL TRANSCRIBERS**



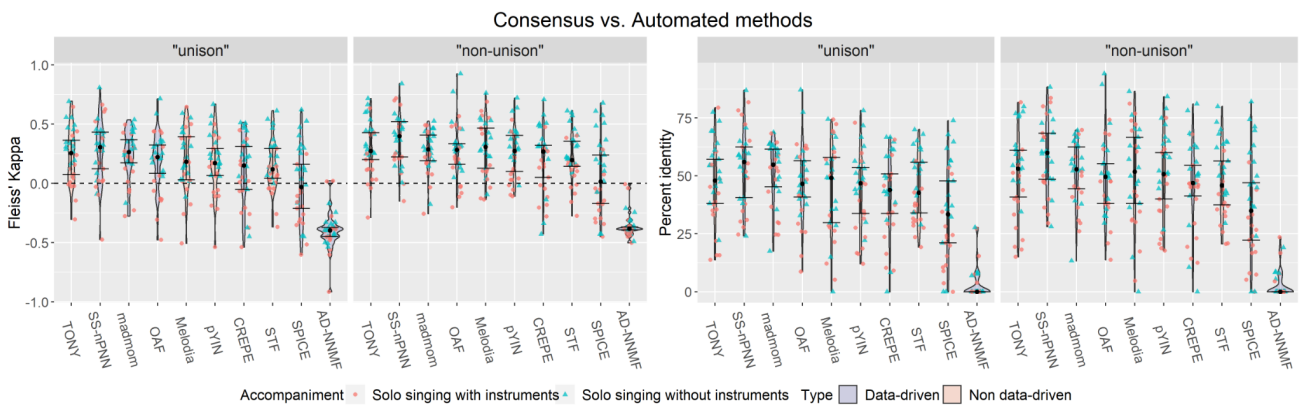
129  
130 **Figure S3.** Pairwise agreement of automated methods vs. human ground-truth transcriptions.



131  
132 **Figure S4.** Pairwise agreement of automated methods vs. human ground-truth transcriptions.

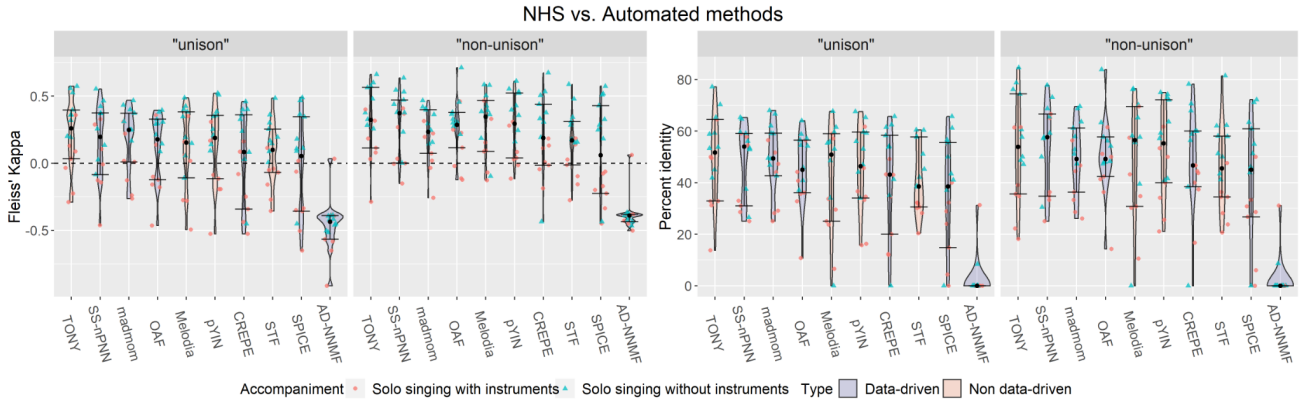


133  
134 **Figure S5.** Pairwise agreement of automated methods vs. human ground-truth transcriptions.



135

136 **Figure S6.** Pairwise agreement of automated methods vs. human ground-truth transcriptions.



137  
138 **Figure S7.** Pairwise agreement of automated methods vs. human ground-truth transcriptions.

139 **L. RESULTS OF HYPOTHESIS TESTING**

140

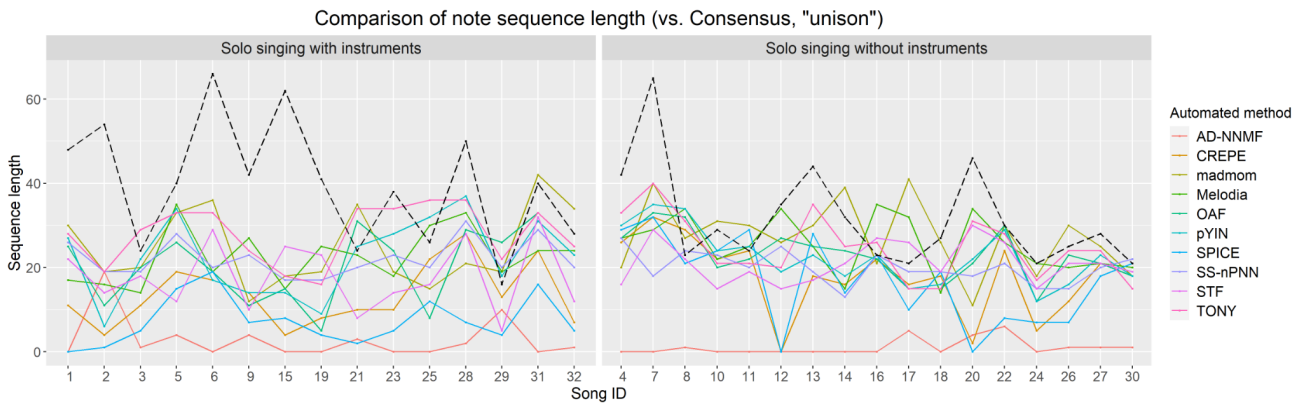
Median of $\kappa = 0$	p-value	$\alpha$ (BH)	ES (g)
Consensus vs. NHS	<b>&lt;0.001</b>	0.010	0.5

141 **Table S2.** Result of the one-sample test.  $\alpha$  (BH) is a threshold adjusted by the Benjamini–Hochberg step-up procedure

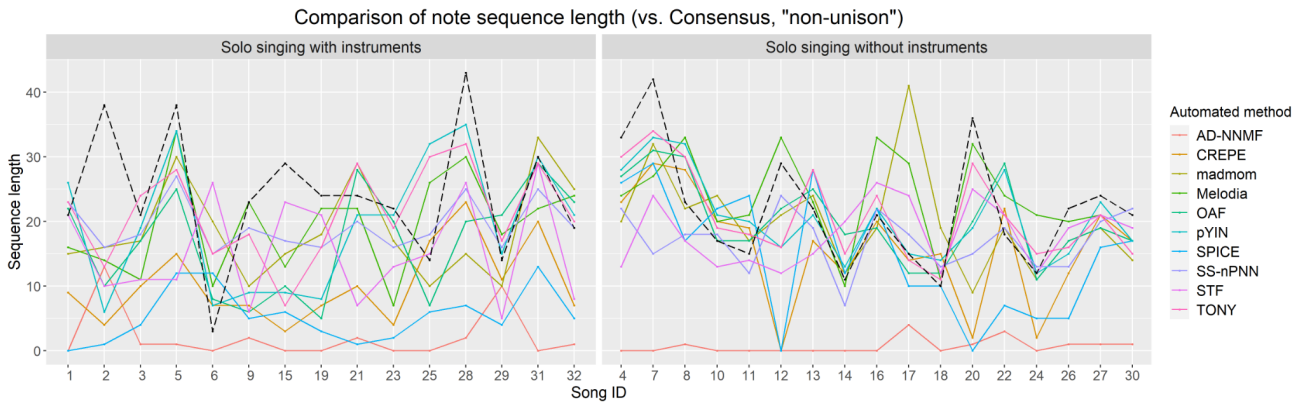
TONY vs.	p-value	$\alpha$ (BH)	$p(H_1 X)$	ES (A)
AD-NMF	<b>&lt;0.001</b>	0.005	<b>1.000</b>	0.985
CREPE	0.104	0.020	0.443	0.623
madmom	0.655	0.040	0.371	0.499
OAF	0.639	0.030	0.152	0.541
SPICE	<b>0.001</b>	0.015	<b>0.970</b>	0.732
SS-nPNN	0.923	0.045	0.145	0.462
Melodia	0.962	0.050	0.193	0.524
STF	0.210	0.025	0.214	0.613
pYIN	0.655	0.035	0.334	0.556

142 **Table S3.** Results of the two-sample tests.  $\alpha$  (BH) is a threshold adjusted by the Benjamini–Hochberg step-up procedure.

143 **M. NOTE LENGTHS OF NOTE SEQUENCES BY AUTOMATED METHODS**

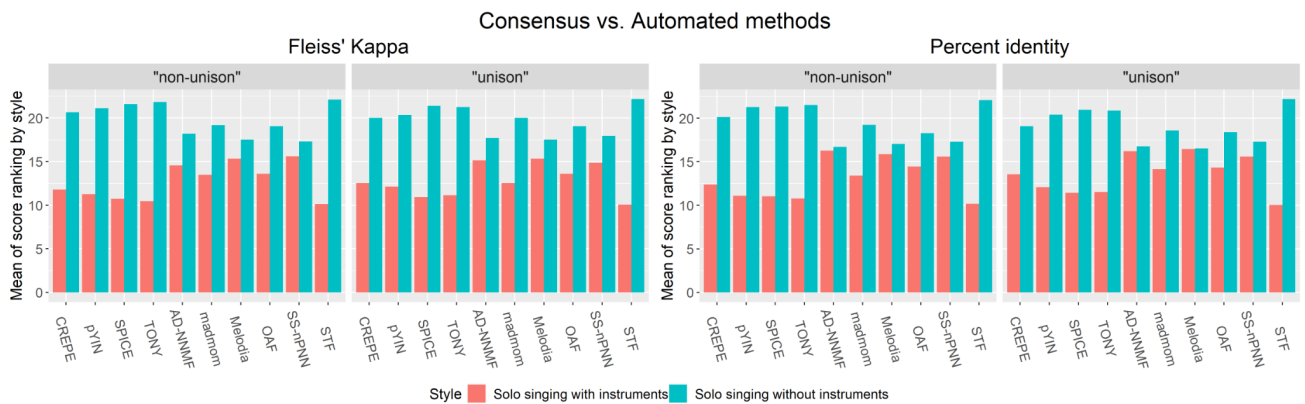


144  
145 **Figure S8.** Lengths of note sequences by automated methods. The dashed line corresponds to the human note sequences,  
146 and the gap against that indicates that notes are segmented more or less than human transcription.



147  
 148 **Figure S9.** Lengths of note sequences by automated methods. The dashed line corresponds to the human note sequences,  
 149 and the gap against that indicates that notes are segmented more or less than human transcription.

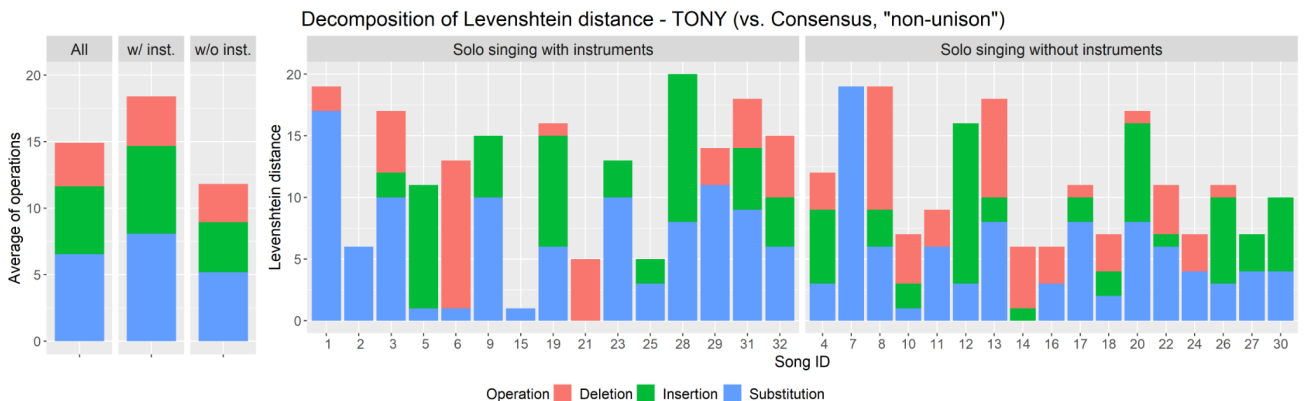
150 **N. DIFFERENCE IN THE ORDER OF AGREEMENT SCORE BY SONG STYLE**



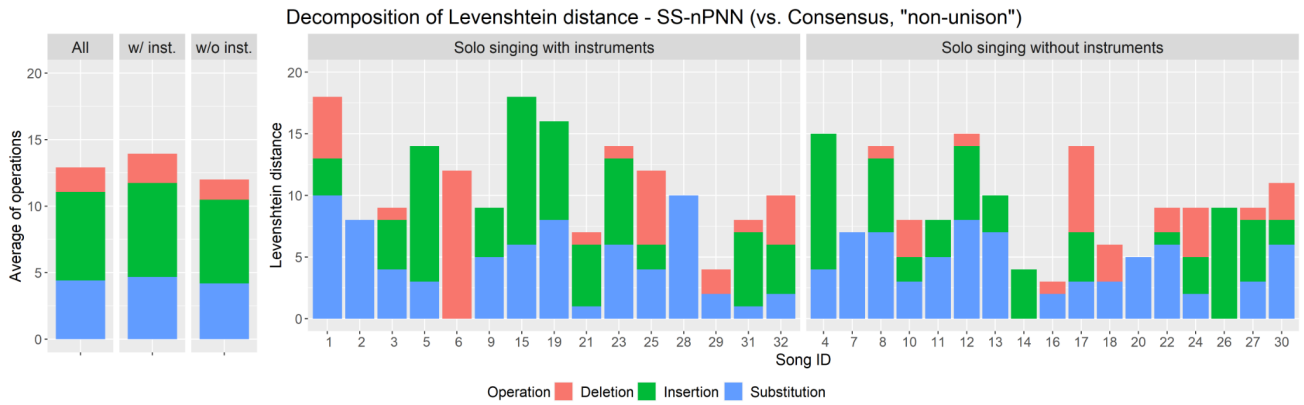
151  
 152 **Figure S10.** Difference of the average of ranking of scores by song styles. Scores of the 32 songs were ranked by de-  
 153 scending order. The gap of average ranking indicates the automated method performed well for one style compared to the  
 154 other.

155 **O. FACTORS AND PATTERNS OF DISAGREEMENT BY LEVENSHTAIN DISTANCE**

156 The below figures show varying patterns of disagreement among the note sequences of human and automated methods.  
 157 We picked up 4 automated methods as representative samples. Furthermore, we chose the "non-unison" version to be  
 158 able to evaluate the F0 prediction performance more directly.



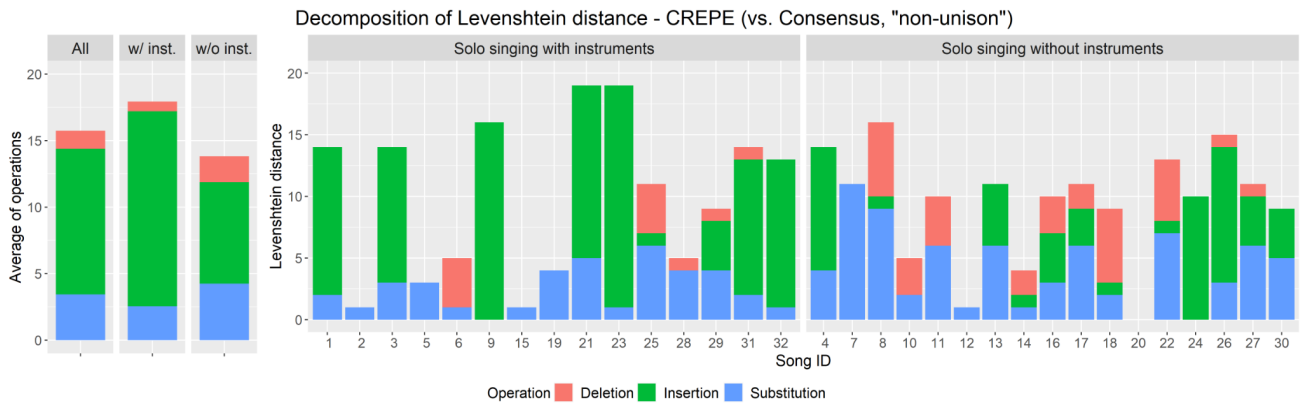
159  
 160 **Figure S11.** Type of disagreement decomposed by operation types.



161

**Figure S12.** Type of disagreement decomposed by operation types.

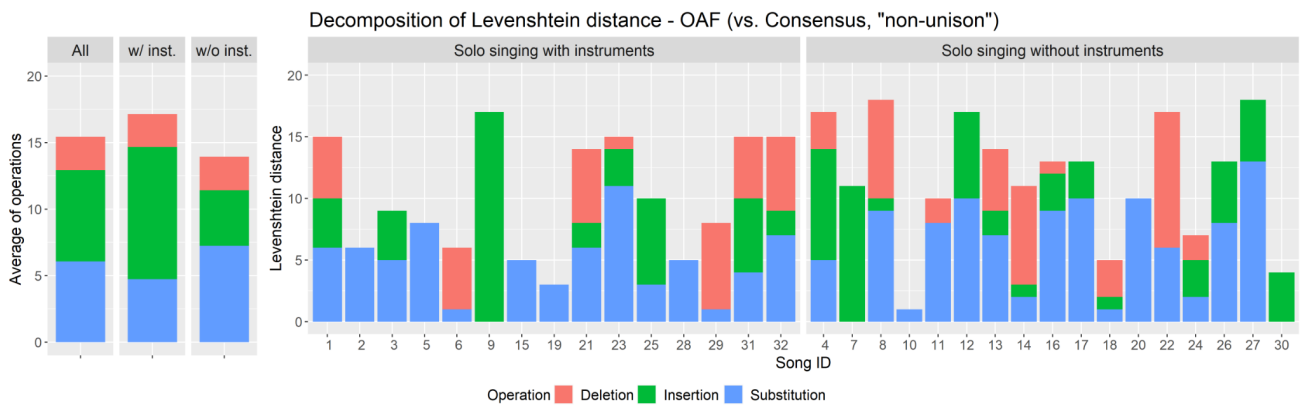
162



163

**Figure S13.** Type of disagreement decomposed by operation types.

164



165

**Figure S14.** Type of disagreement decomposed by operation types.

166

167 **P. SUMMARY OF TOP 10 OVERLAPPED BEST AGREEMENT RESULTS BY AUTOMATED METHODS**

168 We first picked up the top 10 agreement songs in reference to our consensus note sequences from each automated method.

169 After that, we further picked up the top 10 overlapping songs from that result.

Song	Song style	# of top-10 ranking in	Max $\kappa$	Automated method
NAIV-054	Solo singing without instruments	8	0.59	Melodia
NAIV-117	Solo singing without instruments	8	0.81	SS-nPNN
T5468R28	Solo singing without instruments	8	0.56	TONY
T5522R80	Solo singing without instruments	8	0.72	OAF
T5528R18	Solo singing with instruments	8	0.62	SS-nPNN
NAIV-021	Solo singing without instruments	7	0.56	TONY
NAIV-029	Solo singing with instruments	5	0.65	TONY

T5482R03	Solo singing with instruments	5	0.40	TONY
NAIV-075	Solo singing without instruments	4	0.47	madmom
T5421R17	Solo singing with instruments	4	0.67	SS-nPNN

170

**Table S4.** Results by the “unison” note sequence version.

Song	Song style	# of top-10 ranking in	Max $\kappa$	Automated method
NAIV-054	Solo singing without instruments	9	0.93	OAF
NAIV-104	Solo singing without instruments	8	0.58	CREPE
NAIV-117	Solo singing without instruments	8	0.84	SS-nPNN
T5468R28	Solo singing without instruments	8	0.67	TONY
T5522R80	Solo singing without instruments	7	0.77	OAF
T5528R18	Solo singing with instruments	7	0.70	SS-nPNN
NAIV-021	Solo singing without instruments	6	0.61	pYIN
NAIV-029	Solo singing with instruments	4	0.64	TONY
T5421R17	Solo singing with instruments	4	0.67	SS-nPNN
T5487R13	Solo singing with instruments	4	0.72	SS-nPNN

171

**Table S5.** Results by the “non-unison” note sequence version.172 **Q. REFERENCES**

- 173 [50] I. Djurovic, and L. J. Stankovic, “An algorithm for the Wigner distribution based instantaneous frequency estimation  
174 in a high noise environment,” *Signal Processing*, vol. 84, no. 3, pp. 631–643, 2004, doi: 10.1016/j.sigpro.2003.12.006.
- 175 [51] D. Gusfield, "Efficient methods for multiple sequence alignment with guaranteed error bounds," *Bulletin of Mathe-*  
176 *matical Biology*, vol. 55, no. 1, pp.141-154, 1993.
- 177 [52] G. Yona, *Introduction to Computational Proteomics*, Boca Raton, FL, USA: Chapman and Hall/CRC, 2010.
- 178 [53] M. L. McHugh, "Interrater reliability: the kappa statistic," *Biochemia Medica*, vol. 22, no. 3, pp. 276-282, 2012.
- 179 [54] W. R. Pearson, "An Introduction to Sequence Similarity (“Homology”) Searching," *Current Protocols in Bioinform-*  
180 *atics*, vol. 42, no. 1, pp. 3.1.1-3.1.8, 2013, doi: 10.1002/0471250953.bi0301s42.
- 181 [55] A. N. Pettitt, "A Two-Sample Anderson-Darling Rank Statistic," *Biometrika*, vol. 63, no. 1, pp. 161-168, 1976.
- 182 [56] C. C. Holmes, F. Caron, J. E. Griffin, and D. A. Stephens, “Two-Sample Bayesian Nonparametric Hypothesis Testing,”  
183 *Bayesian Analysis*, vol. 10, no. 2, pp. 297–320, 2015.
- 184 [57] J. Cohen, *Statistical Power Analysis for the Behavioral Sciences, 2nd Ed.*, New York, NY, USA: Routledge, 1988.
- 185 [58] J. Ruscio, "A Probability-Based Measure of Effect Size: Robustness to Base Rates and Other Factors," *Psychological*  
186 *Methods*, vol. 13, no. 1, pp. 19–30, 2008, doi: 10.1037/1082-989X.13.1.19.
- 187 [59] J. Cohen, “A power primer,” *Psychological Bulletin*, vol. 112, no. 1, pp. 155-159, 1992.
- 188 [60] Y. Benjamini, and Y. Hochberg, "Controlling the False Discovery Rate: A Practical and Powerful Approach to Mul-  
189 tiple Testing," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 57, no. 1, pp. 289-300, 1995.

190



## RELIABILITY OF AUTOMATED AND HUMAN TRANSCRIPTIONS OF GLOBAL SONGS

Embargoed registration



## Overview

Files

Wiki

Components 0

Links 0

Analytics

Comments 0



## Study Information



## Hypotheses

We will examine the following two questions.

- Q1) To what degree do human transcriptions agree?  
 Q2) Which automated method agrees to the greatest extent with the transcription of non-Western music by humans?

Regarding the first question, based on our pilot study, we hypothesize that inter-rater reliability among human's transcriptions would be significantly larger than 0. We will confirm it by performing the one-tailed Sign Test assuming the median of Fleiss' Kappa to be 0 (i.e. no agreement) as a null hypothesis. Fleiss' Kappa will be computed as the agreement between our collaborator's consensus note sequences and the note sequences created from Natural History of Song's consensus transcription for each song. Transcriptions of 16 songs from the Natural History of Songs dataset will be used to create note sequences. Since transcription would certainly involve subjective interpretation, we would like to confirm its reliability by comparing two transcriptions created by different groups of experts.

Regarding the second question, we assume TONY (Mauch et al., 2015) would be the most competitive method. Therefore, we will test whether the Fleiss' Kappa by TONY is superior to other methods by performing a two-sample test with each of the rest of the methods and TONY. We will compare our consensus note sequences and note sequences created by each automated method to compute Kappa per song. We have two reasons for choosing TONY as a baseline. One reason is based on the empirical observation from our pilot study that TONY achieved the highest median of Kappa with our consensus transcription. Secondly, TONY is the only note-level transcription method designed for vocal melody amongst the selected methods, and our transcription analysis agreement analysis will be using human vocal melody samples and will be undertaken at note-level, not frame-level. Hence, we expect TONY is the most well-designed method as a baseline.

## Design Plan

## Study type

Observational Study - Data is collected from study subjects that are not randomly assigned to a treatment. This includes surveys, "natural experiments," and regression discontinuity designs.

## Blinding

Personnel who analyze the data collected from the study are not aware of the treatment applied to any given group.

## Is there any additional blinding in this study?

Human transcriptions were collected from three Japanese collaborators who are all receiving/received professional musical training. From the blinding protocol viewpoint, they received the explanation of the purpose of research along with the instruction of transcription tasks, but they are not aware how that their transcriptions will be compared and will be evaluated, so the transcribers are not rigorously blinded, but we consider there would be no clear causality which affects the accuracy and quality of their transcriptions. They do not have any information regarding music to be transcribed other than provided audio wav files.

## Study design

We will perform the following patterns of agreement evaluation. We compare note sequences as explained below for each evaluation. However, we will use a subset of those for our hypothesis testing.

- Q1) To what degree do human transcriptions agree?  
 1) Pairwise agreement of note sequences by two of the three collaborators.  
 2) Pairwise agreement between note sequences of our consensus transcription and the three collaborators' transcriptions (i.e. consensus vs. individual transcription).  
 Q2) Which automated method agrees to the greatest extent with the transcription of non-Western music by humans?  
 3) Pairwise agreement between note sequences of human transcription (our collaborators' transcription plus our consensus transcription) and machine transcription.

## Contributors

## Description

Obtaining musically interpretable representation of sound such as score notation to conduct cross-cultural analysis of music is essential to explore its diversity and universality. However, transcription is usually conducted manually by ear, which is time-consuming and subjective. In this study, we examine how much automatic music transcription/vocal melody extraction systems can address this gap with song excerpts from various cultures. To do so, we conduct a series of experiments evaluating the degree of agreement, not only among human transcribers, but also between human transcribers and 10 automatic music transcription/vocal melody extraction methods. Pre-trained models are used for machine learning-based methods. We use 32 monophonic/homophonic melodies in the Natural History of Song dataset (Mehr et al., 2019) and the Cantometrics songs dataset equally sampled from each region (Asia, Middle East, North America, etc.). Both human and machine transcribers only rely on the audio files to create transcriptions. Equal temperament is used to discretize frequency as notes although its appropriateness is not definitive. We use Fleiss' Kappa, percent melodic identity (May, 2004; Savage and Atkinson, 2015) and Levenshtein distance to quantify the agreement, but we will mainly use Fleiss' Kappa for hypothesis testing.

Show less

## Registration type

OSF Preregistration

## Date registered

April 20, 2021

## Date created

April 20, 2021

## Registered from

osf.io/vc8g7

## Category

Analysis

## Publication DOI

No publication DOI

## Subjects

Engineering  
 Social and Behavioral Sciences  
 Music Arts and Humanities  
 Computational Engineering

192

193

194

195

The following evaluations are only applicable to songs from Natural History of Song (NHS) dataset.

- 4) Pairwise agreement between note sequences of NHS consensus transcription and the three collaborators' transcriptions.
- 5) Pairwise agreement between note sequences of NHS consensus transcription and our consensus transcription.
- 6) Pairwise agreement between note sequences of NHS consensus transcription and machine transcription.

We will create the sequence of note from each transcription of a single vocal melody by discarding the note value and rest information, and ordering the pitch class with octave information by the order of appearance. All transcriptions will use twelve-tone equal temperament to discretize the frequency although its appropriateness will be addressed in the paper. In addition, twelve-tone equal temperament has enharmonic equivalent pitch classes, so we will only use flat notes for the same sounding sharp and flat notes to create note sequences.

We created a consensus note sequence of each song by the following steps. Firstly, we automatically aligned the note sequences and then performed manual correction with rhythmic information. We will explore both, but our prediction will be based on the consensus transcription created by manual alignment. Secondly, disagreements of each note among note sequences were resolved by majority rule. If there is a note that is different in all three note sequences, we asked our collaborators via email to choose which note would fit the consensus notes selected by the majority rule. If still disagreement remained, we chose the median of the pitch from the disagreeing notes. However, subjective decisions were made when the alignment gap involved in the disagreement. We also confirmed the soundness of the resultant consensus transcriptions whose disagreement was resolved by their choice and our subjective decision. We further updated some transcriptions when their change proposal based on their soundness evaluation was regarded as more appropriate.

The degree of agreement among note sequences will be quantified by Fleiss' Kappa, percent identity (May, 2004; Savage and Atkinson, 2015), and Levenshtein distance since the notion of agreement between two sequences can be interpreted in several ways. However, we will only use Fleiss' Kappa for our hypothesis testing.

In order to calculate Kappa and percent identity, we need to align sequences to be evaluated. Therefore, we will perform pairwise alignment to create the alignment of note sequences by the Needleman-Wunsch algorithm using 0.0 for gap opening penalty, -1.0 for gap extension penalty and -1.0 for mismatch (substitution) penalty. This is a linear gap setting, and the alignment score is equivalent to Levenshtein distance whose operations (i.e. insertion, deletion, substitution) are all equally weighted. We will use octave information for the evaluation, so the element of the sequence will consist of two characters: pitch class and octave level (e.g. "A4"). When multiple sequence alignment is necessary for creating the baseline of consensus note sequences, we use the center star method to solve alignment heuristically since the computation of the global optimal multiple sequence alignment is not feasible due to its computational complexity (Gusfield, 1993; Yona, 2010). The center sequence is determined by the sum-of-pairs scoring (Gusfield, 1993; Yona, 2010), and each score is calculated by the Needleman-Wunsch algorithm as described above.

We selected 10 automatic music transcription / vocal melody extraction / pitch detection methods in total. Considering the difference in the approach of the pitch estimation, our selection consists of automated methods from non-data-driven models and data-driven models. In this study, we call methods data-driven if the model employs a machine learning method (such as artificial neural networks) to learn model parameters from data in a training step. On the other hand, we call the others non-data-driven. Table 1 (attached image file) summarizes the selected automated methods. Regarding TONY, we use frame-level estimation of TONY as pYIN (Mauch and Dixon, 2014) and the note-level estimation as TONY.

- Regarding the note sequence by automated methods, we will apply the following post-processing procedures against the output of each method to smoothen the pitch contour.
- Regarding the methods which do not quantize F0 to twelve-tone equal temperament, the estimated F0 is rounded to the nearest frequency of the twelve-tone equal temperament.
- Regarding the methods which do not estimate note duration or note tracking, a median filter with the length of 0.25 seconds is applied to smooth the pitch contour. Furthermore, the sequences of F0 shorter than 0.15 seconds are ignored from transcription. 0.15 is determined to make the length of note sequence similar to the humans' sequences. These parameters were tuned to minimize the possibility that the automated methods would produce long sequences made up of unrealistically short notes as a by-product of the instability of pitch targets in human singing. If the unit of the discrete time interval of generated time-frequency representation is less than 0.01 second, decimation is applied to make the interval close to 0.01 second to smoothen the pitch contour.
- OAF estimates onset and offset of note, but it is fairly precise, so the above post-processing is applied to make a more meaningful comparison with the other methods.
- Regarding the methods predicting multiple pitches in a single timeframe, we apply the following steps to obtain the stream of single pitch prediction. Firstly, we observe that these methods tend to predict an overtone as a separate note, so the frequency range of the melody is manually specified, and the F0 prediction out of this range was removed. After that, the Viterbi algorithm is applied to the remaining multi-pitch F0 prediction results to obtain the dominant time-frequency energy sequence as a melody (Djurovic and Stankovic, 2004).
- Regarding CREPE, F0s having a confidence score larger than or equal to 0.8 are picked up.
- We use a song excerpt as the input of automated methods to obtain the pitch estimation of a specified 14-second segment. However, pitch estimation process would depend on the information available on the broader time range of audio data to estimate the F0 of local time-frame, so feeding an entire song as input and extracting the target segment from its output will produce different pitch estimation results. In this study, we only have the excerpt of songs regarding the NHS, so we decided to consolidate the input by an audio excerpt.

#### Affiliated institutions

This registration has no affiliated institutions

#### License

CC-BY Attribution 4.0 International

#### Tags

No tags

#### Citation

osf.io/bjemd

Furthermore, since most methods are not designed to detect the segmentation of voicing, the selected automated methods do not segment unison intervals as like human transcribers. Therefore, we will also create a transcription which discards repeated notes and treats the notes of the unison interval as a tied single note (i.e. "C3C3F3G3G3C3" becomes "C3F3G3C3"). We call this version note sequence without the unison interval, and the original version as note sequence with unison interval. This treatment enables us to evaluate how much the pitch estimation itself, which is a baseline function of automatic transcription, determines performance. We will explore both but will conduct hypothesis testing using the note sequences with the unison interval version.

- [AMT\\_table.png](#)

#### Randomization

This study is not a randomized study.

## Sampling Plan

#### Existing Data

Registration prior to analysis of the data

#### Explanation of existing data

We will use two datasets to sample songs from various societies: Natural History of Song (NHS) and Cantometrics (<http://theglobaljukebox.org>). We judged the appropriateness of the datasets by noting that these two datasets cover traditional vocal music at a very broad range of regions. However, the NHS only provides 14-second excerpts of songs, so we also manually extract 14-second segments from randomly sampled Cantometrics songs to consolidate the song length.

#### Data collection procedures

To examine the reliability of transcription of non-Western music by both humans and computers, 16 songs were sampled from the publicly available 14-second excerpts of the NHS and manually extracted 14-second excerpts of the Cantometrics audio files, respectively. Sampling was randomly conducted using the following criteria.

- 1) NHS categorizes the globe into the eight regions (North America, Asia, etc.). Songs are sampled equally from each region (i.e. 4 songs per region).
- 2) In order to assess the capability of note listening in a broad situation, songs are sampled to consist of solo singing without instruments and solo singing with instruments.

As a result, we sampled two songs for each category (solo a cappella, or solo with instrumental accompaniment) from 8 regions. However, the NHS dataset contained no audio recordings of solo singing with instruments in the Middle East region, so two solo singing without instrument examples were chosen from this region instead.

*No files selected*

#### Sample size

32 songs in total (16 songs without instruments and 16 songs with instruments).

#### Sample size rationale

As explained below, we assume Kappa could be non-normally distributed so our testing method is nonparametric and it is challenging to assess the power. It would be desirable to collect more samples with further variation criteria such as level of noise and singing style to thoroughly evaluate the performance of automated methods. However, there is no previous research evaluating the performance of automated methods by measuring the agreement between human transcriptions with non-Western music, so we would like to set these samples as the tentative baseline. There is also time-constraint for the amount of transcriptions we can ask our collaborators to create.

#### Stopping rule

NA

## Variables

#### Manipulated variables

NA

*No files selected*

#### Measured variables

We will compute the agreement among transcriptions by Fleiss' Kappa, percent identity and Levenshtein distance as variables we will measure to investigate our questions.

*No files selected*

#### Indices

NA

*No files selected*

## Analysis Plan

### Statistical models

The underlying distribution of inter-rater reliability coefficients is considered to depend on the raters (i.e. transcribers) and subjects (audio recording) (Gwet, 2015).

Furthermore, our agreement metrics are collected from various combinations of transcribers and audio samples, and the domain of Kappa coefficients is finitely bounded, so the resultant distributions of agreement metrics would not necessarily fit normal or location-scale family distributions.

Based on the above assumption, we consider the appropriate testing methods to handle the metrics to be nonparametric methods. We choose the Sign Test for one-sample test scenario, and the two-sample Anderson-Darling test (Pettitt, 1976) and two-sample Bayesian nonparametric testing using Pólya trees (Holmes et al., 2015) for two-sample test scenarios. Regarding the two-sample test, we assess the probability of type I error by the two-sample Anderson-Darling test. Besides, to complement the lack of information about how much we can be confident in accepting alternative hypotheses, we also employ Bayesian hypothesis testing. Although these two tests are different procedures, both are proved to be asymptotically consistent under the null hypothesis ( $F(x)=G(x)$ ) and the alternative hypothesis ( $F(x) \neq G(x)$ ) (Holmes et al., 2015; Pettitt, 1976). We set  $c=1$  and the normal distribution as the centering distribution as the parameters of the Pólya trees. However, we use the mean and standard deviation to create partitions of samples instead of the median and quantiles used in the original setting (Holmes et al., 2015). We set the equal probability for the null hypothesis and the alternative hypothesis ( $= 0.5$ ) as the prior distribution of our Bayesian hypothesis testing.

Regarding the effect size to be used for our nonparametric two-sample test, we choose the probability-based effect size measure  $A$  which is known as the probability of one group's superiority over another (Ruscio, 2008). The probability-based effect size uses empirical distributions of data to quantify how much data in a group takes a larger value than another group, and it is robust to violations of the parametric assumptions. Note that  $A$  can be converted to a common standardized mean difference such as Cohen's  $d$  if the normality assumption of data holds.

Regarding the first question Q1, we will conduct the one-tailed one-sample Sign Test as explained in the Hypothesis section.

Regarding the second question Q2, firstly, we will confirm whether TONY could achieve the highest median with consensus note sequences compared to the other methods. We will only report the results of Kappa since this metric would be most relevant to the concept of reliability. Each automated method will have 32 songs (16 Cantometrics and 16 NHS) of Fleiss' Kappa calculated with consensus transcription, so we will measure the median of Kappa of each automated method over 32 songs. After that, we will test whether Fleiss' Kappa of TONY has a statistically significant difference from the others using the two-sample test methods described above. We will also check  $A$  to assess whether the data points of the empirical distribution of the best-performing method's Kappa are generally larger than the other ones.

*No files selected*

### Transformations

NA

### Inference criteria

There will be 10 null hypothesis significance tests in our analysis: one-sample Sign Test  $\times 1 +$  two-sample Anderson-Darling test  $\times 9$  (machine pairs). Since our discussion on the reliability of transcription will be interrelated to these test results, we will use the False Discovery Rate method to control the p-value threshold for all hypothesis tests regarding these as multiple testing and simultaneous inference. In particular, we will use the Benjamini-Hochberg step-up procedure at level  $\alpha = 0.05$  as the threshold to determine the rejection of 10 null hypotheses.

We will interpret the Bayesian test at least substantially supports the alternative hypothesis if the posterior probability exceeds 0.8 which corresponds to the Bayes factor  $= 4$  in our setting (i.e. the prior distribution being equally weighed to the null and alternative hypothesis). We also check the effect size to grasp the superiority of agreement and whether it does not contradict with p-value and the posterior probability.

### Data exclusion

As long as transcriptions are created for the entire section of songs for both human and machine cases, we will use all data for our analysis.

### Missing data

NA

### Exploratory analysis

We will analyze the difference between note sequences by decomposing it with the operation of Levenshtein distance: insertion, deletion, and substitution. This would help us systematically explore if there are specific patterns of difference among transcriptions.

We will also explore the following hypothesis where appropriate.

- Is there a difference in the degree of the agreement between data-driven methods and non-data-driven methods?
- Do humans' transcriptions differ more for short-length notes compared to longer-length notes?
- Do humans' transcriptions differ more for rarely appearing compared to more commonly appearing notes?

## Other

### Other

Bock, S., Korzeniowski, F., Schluter, J., Krebs, F., & Widmer, G. (2016). Madmom: A new Python Audio and Music Signal Processing Library. Proc. 24th ACM International Conference on Multimedia, Amsterdam, Netherlands.

Cheng, T., Mauch, M., Benetos, E. and Dixon, S. (2016). An attack/decay model for piano transcription. Proc. 17th International Society for Music Information Retrieval Conference (pp. 584-590), New York, NY, USA.

Djurovic, I. & Stankovic, L. J. (2004). An algorithm for the Wigner distribution based instantaneous frequency estimation in a high noise environment. *Signal Processing*, 84(3), 631-643, <https://doi.org/10.1016/j.sigpro.2003.12.006>.

Gfeller, B., Frank, C., Roblek, D., Sharifi, M., Tagliasacchi, M. & Velimirović, M. (2020). SPICE: Self-Supervised Pitch Estimation. *IEEE/ACM Transactions on Audio, Speech, and Language*, 28, 1118-1128. <https://doi.org/10.1109/TASLP.2020.2982285>.

Gusfield, D. (1993). Efficient methods for multiple sequence alignment with guaranteed error bounds. *Bulletin of Mathematical Biology*, 55(1), 141-154.

Gwet, K. L. (2015). *Handbook of Inter-Rater Reliability. The Definitive Guide to Measuring the Extent of Agreement Among Raters*, 4th edition. Gaithersburg, MD, USA: Advanced Analytics, LLC.

Hawthorne, C., Elsen, E., Song, J., Roberts, A., Simon, I., Raffel, C., Engel, J., Oore, S. & Eck, D. (2018). Onsets and Frames: Dual-Objective Piano Transcription. Proc. 19th International Society for Music Information Retrieval Conference (pp. 50-57), Paris, France.

Holmes, C. C., Caron, F., Griffin, J. E. and Stephens, D. A. (2015). Two-Sample Bayesian Nonparametric Hypothesis Testing. *Bayesian Analysis*, 10(2), 297-320.

Kim, J. W., Salamon, J., Li P. & Bello, J. P. (2018). Crepe: A Convolutional Representation for Pitch Estimation. Proc. 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (pp. 161-165), Calgary, AB, Canada. <https://doi.org/10.1109/ICASSP.2018.8461329>.

Lu, W.-T. & Su, L. (2018). Vocal melody extraction with semantic segmentation and audio-symbolic domain transfer learning. Proc. 19th International Society for Music Information Retrieval Conference (pp. 521-528), Paris, France.

Mauch, M., Cannam, C., Bittner, R., Fazekas, G., Salamon, J., Dai, J., Bello, J. & Dixon, S. (2015). Computer-aided melody note transcription using the tony software: Accuracy and efficiency. Proc. 1st International Conference on Technologies for Music Notation and Representation, Paris, France.

Mauch, M. & Dixon, S. (2014). PYIN: A Fundamental Frequency Estimator using Probabilistic Threshold Distributions. Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (pp. 659-663), Florence, Italy.

May, A. C. W. (2004). Percent Sequence Identity: The Need to Be Explicit. *Structure*, 12, 737-738. <https://doi.org/10.1016/j.str.2004.04.001>.

Mehr, S. A., Singh, M., Knox, D., Ketter, D. M., Pickens-Jones, D., Atwood, S., Lucas, C., Jacoby, N., Egner, A. A., Hopkins, E. J., Howard, R. M., Hartshorne, J. K., Jennings, M. V., Simson, J., Bainbridge, C. M., Pinker, S., O'Donnell, T. J., Krasnow, M. M. & Glowacki, L. (2019). Universality and diversity in human song. *Science*, 366, eaax0868. <http://dx.doi.org/10.1126/science.aax0868>.

Pettitt, A. N. (1976). A Two-Sample Anderson-Darling Rank Statistic. *Biometrika*, 63(1), 161-168.

Ruscio, J. (2008). A Probability-Based Measure of Effect Size: Robustness to Base Rates and Other Factors. *Psychological Methods*, 13(1), 19-30. <https://doi.org/10.1037/1082-989x.13.1.19>.

Salamon, J. & Gomez, E. (2012). Melody Extraction from Polyphonic Music Signals Using Pitch Contour Characteristics. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(6), 1759-1770. <https://doi.org/10.1109/TASL.2012.2188515>.

Savage, P. E. & Atkinson, Q. D. (2015). Automatic Tune Family Identification by Musical Sequence Alignment. Proc. 16th International Society for Music Information Retrieval Conference (pp. 162-168), Málaga, Spain.

Su, L. & Yang, Y. (2015). Combining Spectral and Temporal Representations for Multipitch Estimation of Polyphonic Music. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(10), 1600-1612. <https://doi.org/10.1109/TASLP.2015.2442411>.

Yona, G. (2010). *Introduction to Computational Proteomics*, Boca Raton, FL, USA: Chapman and Hall/CRC.



This file is part of a registration and is being shown in its archived version (and cannot be altered). The [active file](#) is viewable from within the [live project](#).

This registration is currently embargoed. It will remain private until its embargo end date, Sunday, Jul 11, 2021.

## AMT\_table.png (Version: 1)

Check out Delete Download View Revisions

Filter

- RELIABILITY OF AUTOMATED AN...
- OSF Storage (United States)
  - Archive of OSF Storage
    - Figure and Table
      - AMT\_table.png

Tags

Add a tag to enhance discoverability

Method	Target sound	Unit	Category
pYIN (Mauch and Dixon, 2014)	Monophonic vocal	Frame	Non data-driven: parameters specified manually
TONY (Mauch et al., 2015)	Monophonic vocal	Note	
Melodia (Salamon and Gomez, 2012)	Vocal melody	Frame	
STF (Su and Yang, 2015)	Multiple 12-tone ET	Frame	
CREPE (Kim et al., 2018)	Monophonic vocal	Frame	Data-driven: parameters optimized by training with datasets.
SPICE (Gfeller et al., 2020)	Monophonic vocal	Frame	
SS-nPNN <sup>1</sup> (Lu and Su, 2018)	Vocal melody	Frame	
AD-NMF <sup>2</sup> (Cheng et al., 2016)	Multiple piano sound	Note	
OAF <sup>3</sup> (Hawthorne et al., 2018)	Multiple piano sound	Note	
madmom (Bock et al., 2016)	Multiple piano sound	Note	

**Table 1.** Summary of the selected automated methods. Unit indicates if the F0 estimation is frame-level or note-level that the latter predicts onset and offset timing.

