Variation in Transcribing Heritage Cantonese







Dependent variables	Homeland change?	Heritage change?
(I): /i/-/ɪ/ split (Tse 2019)	Yes	Yes
(Y): /y/-/u/ merger (Tse 2019)	No	Yes
(E): /ε/ split by coda (Tse 2019)	No	No
(O): /ɔ/ split by coda (Tse 2019)	Yes	No
(VOT) (Tan & Nagy 2017)	No	No
(N): $/n/ \rightarrow /l/$ (Nagy in press)	Yes	No
(CL): specializing to singular nouns (Nagy & Lo 2019)	No	Yes
(PRODROP): less pro-drop (Nagy et al. 2011, Nagy in press)	No	No
(MOTION): satellite-/verb-framing (Leung 2022)	No	Yes

WC





F (Heritage	-			
(Toronto)	Gen 1	11	7	18
	(10101110)	Gen 2	10	15	25
		Gen 3	2	4	6
F (Homeland Hong Kong	g)	11	3	14
Т	Fotal		34	29	63



Methods

- 1) Corpus-based imputation
- 2) Machine-transliteration (Diep 2022)

9

WOC 2023







Machine-Transliteration Data

We train on publicly-available corpora of Cantonese with both character and jyutping transcriptions, such as:

- Hong Kong Cantonese Corpus (Luke & Wong 2015)
- Child Heritage Chinese Corpus (Mai & Yip 2017)
- Guthrie Bilingual Corpus (Guthrie 1983)
- HKU-70 Corpus (Weizman & Fletcher 2000)
- Lee-Wong-Leung Corpus (Lee et al. 1991–94)
- Leo Corpus (Mai & Yip 2022)
- Paidologos Corpus: Cantonese (Edwards & Beckerman 2008)
- Yip-Matthews Bilingual Corpus (Yip & Matthews 2007)

Over 4 million Chinese words from over 1,232 interviews

WOC 2023

13











CER for Transliteration Methods

We report similar average CER (character error rate) for each transliteration method across homeland and heritage speakers

	Corpus-based Imputation	Machine Transliteration
Homeland	0.30	0.21
Heritage	0.31	0.26
OC 2023		





References

Census and Statistics Department. (2014). Hong Kong monthly digest of statistics: Use of language in Hong Kong in 2012. http://www.statistics.gov.hk/pub/B71406FB2014XXXXB0100.pdf. Accessed April 15, 2017.

Diep, B. (2022). Towards a neural model for jyutping to character transliteration [Course paper for STA497: Readings in Statistics].

Edwards, J., & M. E. Beckman. 2008. Methodological question in studying consonant acquisition. *Clinical Linguistics and Phanetics*, 22(12), 939–958.

Gorman, K., Howell, J., & Wagner, M. (2011). Prosodylab-Aligner: A tool for forced alignment of laboratory speech. *Canadian Acoustics*, 39(3), 192–193.

Guthrie, L. F. (1983). Learning to use a new language: Language functions and use by first grade Chinese-Americans. ARC Associates.

Lee, T. H.T., C. H. Wong, S. Leung, P. Man, A. Cheung, K. Szeto, & C. S. P. Wong. (1991–94). The development of grammatical competence in Cantonese-speaking children. Report of RGC earmarked grant.

Leung, J. R. (2022). Variation in path encoding in motion events in Toronto Heritage Cantonese. University of Pennsylvania Working Papers in Linguistics, 28(2), Article 10.

Luke, K. K., & Wong, M. L. Y. (2015). The Hong Kong Cantonese corpus: Design and uses. Journal of Chinese Linguistics Monograph Series, 25, 312–333.

Mai, Z., & V. Yip. (2017). Acquiring Chinese as a heritage language in English-speaking countries and the Child Heritage Chinese Corpus. International Conference on Bilingualism: Language and Heritage, Chinese University of Hong Kong, Dec. 18.

Mai Z., & V. Yip. (2022). Caretaker input and trilingual development of Mandarin, Cantonese and English in early childhood (1;6-2;11). International Journal of Bilingual Education and Bilingualism, 25(9), 3389–3403.

Nagy, N. (2017). Heritage language speakers in the university classroom, doing research. In P. Trifonas & T. Aravossitas, eds. International handbook on research and practice in heritage language education. Springer.

Nagy, N. (in press). Heritage languages: Extending variationist approaches. Cambridge University Press.

WOC 2023

23

