

# Patterns of variation in Jakarta Indonesian: Linguistic and social dimensions

Abigail C. COHN<sup>#</sup>, Rachel C. VOGEL<sup>#</sup> & Maya Ravindranath ABTAHIAN<sup>\*</sup>

<sup>#</sup>Cornell University, <sup>\*</sup>University of Rochester

Colloquial varieties of Indonesian are increasingly becoming the native languages of a significant portion of the Indonesian population. Notable in this regard is Jakarta Indonesian (JI). We seek to examine the nature of variation in this increasingly widely spoken variety based on the Betawi-Jakarta Indonesian corpus (Gil & Tadmor 2014). We investigate variation within a subset of speakers comparing the phonological variables Kurniawan (2018) examined (word-final [a] ~ [e], word-final [h] ~ [ʔ] ~ Ø, and active prefix N- ~ [ŋə]) with the additional variables word-initial [s] ~ Ø and [h] ~ Ø (Cohn & Vogel 2019) and first person singular (1SG) pronouns (Abtahian, Cohn, Djenar & Vogel 2021). Investigation of this new emerging variety demonstrates both inter- and intra-speaker variation for the variables analyzed, but shows that the variables are not all conditioned by the same linguistic and social factors.

## 1. Introduction<sup>1</sup>

Indonesia is a complex multilingual nation with over 700 local and regional languages in use together with the national language Bahasa Indonesia or Indonesian (Eberhard, Simons & Fennig 2022). Recently attention has been paid to the impact of Indonesian on the maintenance and use of local languages and there is clear evidence that these languages small and large are at the risk of endangerment (Musgrave 2014, Ravindranath & Cohn 2014). However, complementing these shifts away from local languages is a very rich and dynamic language ecology of emerging and increasingly important colloquial varieties of Indonesian (see e.g., Manns, Cole & Goebel 2016). In recent decades, these varieties have become the native languages of a significant portion of the population and these shifts are likely to accelerate in the coming years.

There is thus an acute need for documentation of these increasingly widely used varieties. There are certain inherent challenges in this regard, as spoken language is often given less weight or attention than written language by formal institutions and speakers themselves, with spoken language often dismissed as *Bahasa sehari-hari* ‘everyday language’ rather than as a linguistic system in its own right worthy of investigation. This is particularly true in diglossic situations, such as is the case in Indonesia, with the role of formal Indonesian sanctioned and promoted by the state (*Bahasa Indonesia yang baik dan benar* ‘Indonesian that is good and correct’, see Errington 2014, Sneddon 2003a) as the national language, the language of education and the language of government. Some work has been done describing and recognizing varieties of colloquial Indonesian (notably Ewing 2005 and Sneddon 2006), but much remains to be done.

Jakarta Indonesian (JI) is one such colloquial variety, spoken by several million Indonesians in and around the capital Jakarta and increasingly serving as a model for

---

<sup>1</sup> **Acknowledgements** A version of this paper was presented at KOLITA 20. Thanks to the organizers for the invitation and to the audience for their questions and feedback. Special thanks are due to Dr. Ferdinan Kurniawan, as this work very much builds on the methodological and scholarly contributions of his 2018 Cornell dissertation. Thanks also to Dr. Novi Djenar for her work with us on the first person singular. And thanks to two anonymous reviewers and the NUSA editors for their careful reading and feedback.

urban varieties elsewhere in Indonesia as well as a model for youth language (see e.g., Ewing 2019's work on Indonesian as spoken in Bandung, West Java). JI is described as a contact variety between Betawi, the local variety of Malay historically spoken in the Jakarta area, and Standard Indonesian (SI, Ikranagara 1975, Kurniawan 2018, Sneddon 2006, Wouk 1999).

JI is distinguished from Standard Indonesian by a number of lexical, phonological, morphological, syntactic, and discourse features (Ewing 2005, Sneddon 2006). Sneddon describes JI as a social style or register. As is often characteristic of contact varieties, however, the features that distinguish it from Standard Indonesian are described as showing variable realization for speakers. Thus not only is there a need for fuller documentation of the structure of JI, but also of the observed patterns of variation for different facets of the language. Sneddon (2006:1) provides some insight into the patterns of variation for these variables based on a series of recordings, focusing on what he describes as "language as spoken by educated Jakartans in everyday interactions", taking what he calls Colloquial Jakarta Indonesian as "the prestige variety of colloquial Indonesian" and described with reference to SI. Ewing (2005:228) discussing Colloquial Indonesian more generally describes it as "based on the type of language typically used by educated speakers of Indonesian using the language in ethnically (or first-language) mixed, informal interactions. We understand JI in broader terms not limited to educated speakers, needing to be described in its own right.

To extend this work, naturalistic data is critically important. This is because the use of these colloquial features is most likely to be seen in informal casual speech between members of the same language community and they are much less likely to be used in anything perceived to be a more formal setting. Only in such settings can we gain insight into the patterns of variation used by speakers as part of their linguistic repertoire. Furthermore it is widely understood that the seeds of language change are seen in informal casual speech more prominently and earlier than in more formal speech varieties (Labov 2010 *inter alia*). Happily in this regard, the Betawi-Jakarta Indonesian corpus (Gil & Tadmor 2014) provides this sort of data. The corpus, based on data collected between 2004–2012, contains naturalistic conversations involving 143 speakers, recorded in informal settings around Jakarta. The corpus thus offers the opportunity to examine the linguistic and social dimensions of multiple linguistic variables within the same set of JI and Betawi speakers.

Kurniawan (2018) undertakes just this type of documentation and analysis looking at three phonological variables in close detail, based on a subset of 20 individuals characterized as JI speakers, roughly balanced for gender and education. This includes the following three variables:

(1) Kurniawan (2018) three phonological variables

- a. word-final [a] ~ [e]
- b. word-final [h] ~ [ʔ] ~ Ø
- c. active prefix N- ~ [ŋə]

In each case Kurniawan carefully studies the linguistic conditioning of these variables as well as the socio-indexical factors. He focuses in particular on gender, education level, and age, the latter by using two additional corpora, Wallace (1976) for earlier speakers and Gil & Tadmor (2007) for younger speakers.

In our current project, we extend the work done by Kurniawan looking at three additional variables using the BJI corpus and investigating the patterns of usage in the same set of

speakers for these additional variables. Two are phonological alternations involving word-initial consonants (Cohn & Vogel 2019): [s] ~ Ø, taken to be lexicalized in certain grammatical forms (e.g., *saja* ~ *aja* ‘just’, *sampe* ~ *ampe* ‘until’), and [h] ~ Ø, said to be an optional phonological rule of /h/ deletion (e.g., *hari* ~ *ari* ‘day’, *habis* ~ *abis* ‘finished’). The final variable involves different forms of the first person singular (1SG) pronouns (Abtahian, Cohn, Djenar & Vogel 2021), e.g., *saya* ~ *gua/gue* ~ *aku* ~ *kita/kite*).

Looking at these additional variables not only allows us to extend the work done by Kurniawan, but it also enables us to consider the ways that different linguistic variables may or may not tell the same story of their respective patterns of variation and the factors conditioning this variation.

Here, we report on the results of each of these six variables, addressing structural linguistic conditioning, as well as inter- and intra-speaker variation. Broadly, we find that the linguistic variables analyzed are not all conditioned by the same factors. This reflects the complex nature of linguistic variation and highlights the importance of examining a wide range of variables to understand variation within any particular linguistic context. This is of particular significance in the investigation of new emerging varieties such as JI.

The structure of the paper is as follows. In the rest of this introduction we present background on the nature of linguistic variation with examples from the Indonesian context (section 1.1) and Jakarta Indonesian as an emerging colloquial variety (1.2). Then in section 2, we review the findings of Kurniawan (2018) regarding the three variables he studied. In section 3, we present results regarding [h] ~ Ø, [s] ~ Ø and in section 4, we discuss variation in use of 1SG pronouns. In section 5, we compare the variables and discuss the implications for the development of JI and emerging varieties of language more generally.

### 1.1 The nature of linguistic variation

A very common modern linguistic approach to the description and analysis of language focuses on “linguistic competence”, what language users “know” (explicitly or implicitly) about their languages. This view goes back to the early work of Noam Chomsky in *Aspects of the Theory Syntax* (1965:3): “Linguistic theory is concerned primarily with an ideal speaker-listener in a completely homogeneous speech community.” This abstraction or idealization leads to work describing languages such as “English” or “Indonesian”.

The sociolinguistic approach developed and articulated by William Labov starting with his seminal work in the 1960s, building on the earlier work of Dell Hymes, John Gumperz, and others, focuses on interaction of social and linguistic factors in analyzing language structure and use. Labov (2006:380) states that “the linguistic behavior of individuals cannot be understood without knowledge of the communities that they belong to”. To understand language, we need to study both **linguistic** competence and **communicative** competence. We need to be careful about using labels like “English” or “Indonesian”; this leads us to the fraught territory of differentiating between “languages” and “dialects” when a more neutral way to approach this is as “language varieties”.

One of our central goals as linguists is describing observed patterns in particular languages. The patterns might be due to linguistic factors, so for example in the case of phonological patterns, neighboring segments or syllable structure might condition observed alternations.

Take for example the pattern described for Standard Indonesian by Lapoliwa (1981) and others where the word final /k/ in *masuk* ‘enter’ is realized as [ʔ], [masuʔ] rather than [k], as compared to for example, *masuk* [masuk] ‘entry’. We describe this with a

phonological “rule” as shown in (2a):

(2) Realization of word final /k/

- a. (Standard) Indonesian: /k/ → [ʔ] / \_\_ # (Lapoliwa 1981)
- b. some terminology
  - variable: realization of word final k
  - variants: [k, ʔ]
  - factor: linguistic word position, syllable structure

Because we want to consider at the same time linguistic and social conditioning factors we are going to think of these phonological alternations as variables, the way we frame things in sociolinguistic terms. As listed in (2b), we can describe the realization of word final /k/ as a **variable** with the possible realizations or **variants** [k, ʔ] and the relevant **linguistic factor** or phonological conditioning being word position and syllable structure. Notably, there is also variability observed based on the speaker, the situation, and so forth; so there is more going on here than just the phonological conditioning.

The focus of variationist sociolinguistics is understanding the social dimensions, whether determined by the identity of particular speakers, or the choices they might make in terms of context, topic, addressee and so forth. These dimensions of linguistic practice can be studied through the identification of variables that vary based on such factors. Take for example the variable realization of non-nominal negation in JI studied by Sneddon (2006) as shown in (3):

(3) Variable: Realization of non-nominal negation (Sneddon 2006:57)

- a. variants: *tidak*, *enggak*, *kagak*, *ndak*
- b. factors: formal –H(igh) / informal –L(ow)  
also some effects of age

As found in Sneddon’s study, there are four common variants: *tidak*, *enggak*, *kagak*, *ndak*. As shown in Table 1 (repeated from Sneddon 2006: 57, Table 7a), the informal or “L” variant *enggak* is by far the most common in conversations and interviews, while the formal variant *tidak* is almost on par with *enggak* in meetings, which are a much more formal setting. Sneddon also finds some effects of age whereby younger speakers show significant variation in the interview setting.

**Table 1: Sneddon (2006: 57) Table 7a**

	<i>enggak</i>	<i>kagak</i>	<i>ndak</i>	<i>tidak</i>	total	%L
conversations	2207	36	18	48	2309	97.9
interviews	2285	5	1	196	2472	92.7
meetings	160	1	11	145	317	54.3

Many researchers take either a linguistic **or** a sociolinguistic approach. In our project we are interested in considering both in order to better understand how they interact. In order to do this, we need to consider patterns of variation and analyze the linguistic, discourse/pragmatic, and social factors contributing to the patterns. Crucially when we observe patterns of variation, we ask the question of whether this is between speakers or

whether individual speakers show variation, that is, what we term between speaker or “**inter-**” speaker or within speaker or “**intra-**” speaker variation.

Consider an example of **inter**-speaker variation, where patterns of language use differ between speakers speaking the “same” language based on gender, age, ethnic background, educational background, etc. Take for example terms of address.

- (4) Variable: Realization of terms of address
- a. Indonesian: *Kakak* [kakaʔ] ‘older sibling’
  - b. Javanese background: *Mas* [mas] ‘older brother’; *Mbak* [ˈmbaʔ] ‘older sister’
  - c. Batak background: *Ito* [ito] ‘older brother’; *Eda* [eda] ‘older sister’
  - d. factor: ethnic background of speaker

When speaking Indonesian, someone of Javanese background might use *Mas*, *Mbak* and someone of Batak background might use *Ito*, *Eda* indicating the speaker’s ethnic background. Such variation is based on the identity, in this case ethnic background, of the speaker.

We can contrast this with a case of **intra**-speaker variation, where individual speakers speak differently depending on the context, who they are speaking to, the topic being discussed, etc. As noted above there are noticeable differences between “formal” (H) and “informal” (L) ways of speaking in Indonesian, associated with SI and JI respectively. This can be seen in two different ways that the question ‘You’ve bathed, haven’t you?’ might be asked:

(5) Variation in formal and informal speech

- a. ‘You’ve bathed, haven’t you?’
- b. formal: Anda sudah mandi kan?  
2s already bathe neg Q
- c. informal: Lu uda mandi kan?  
2s already bathe neg Q
- d. variable: form of the 2nd person pronoun  
variants: *Anda* (formal or “H”), *Lu* (informal or “L”)
- e. variable: realization of the form (s)udah ‘already’  
variants: *sudah* (formal or “H”), *uda* (informal or “L”)
- f. factor: degree of formality (as determined by situation, addressee, topic)

In our project, we strive to integrate linguistic and sociolinguistic approaches. In order to do this, we need to study variation without making a priori assumptions about the sources of that variation. We need naturalistic data most readily available in the form of corpus data, crucially including speaker metadata (see Cohn & Renwick 2021 for discussion of these desiderata). This also allows us to study language change over time and the relationship between variation and change.

The issue of the nature of variation is all the more interesting in an emerging linguistic variety, such is the case in Jakarta Indonesian (JI). Emerging contact varieties such as JI are likely to have more variation, at least in early stages, due to the fact that speakers may have different primary languages or multilingual repertoires. As the variety becomes spoken more widely as a primary language, we might expect some of this variation to level out. This is why it is particularly interesting to look at such varieties from both a synchronic and diachronic perspective. We turn now to a bit more information about JI

as an emerging colloquial variety.

## 1.2 Jakarta Indonesian as an emerging colloquial variety

Jakarta Indonesian (JI) is a colloquial variety of Malay and part of a complex linguistic landscape in Indonesia, a nation with the fourth largest population in the world (currently estimated to be about 275 million people) and over 700 languages spoken (about 10% of all the languages of the world, Eberhard, Simons & Fennig 2022) across an archipelago of over 14,000 islands. During the second half of the 20th century, Standard Indonesian, a variety of Malay, was developed and promoted as the national language for the new nation-state of Indonesia (Sneddon 2003a). Following a series of extraordinarily successful language planning efforts, Indonesian is now the dominant language for a large and growing percentage of the population (Anwar 1980, Dardjowidjodjo 1998). As of the 2010 census Indonesian was reported to be the second most widely spoken language at home (Ananta et al. 2015), with 42 million speakers, making it the world's 30th most widely spoken native language worldwide. The use of the label "Indonesian" in the census, however, is somewhat complicated by the fact that the term is used to refer to both the standard variety of the language and colloquial varieties associated with different regions. In recent decades, colloquial varieties of Indonesian have become the native languages of a significant portion of the population. These are sometimes, but not always, distinguished by speakers as varieties distinct from Standard Indonesian (Abtahian, Cohn, White & Yanti 2019, Abtahian, Cohn & Yanti 2022) and may be more or less mutually intelligible with SI.

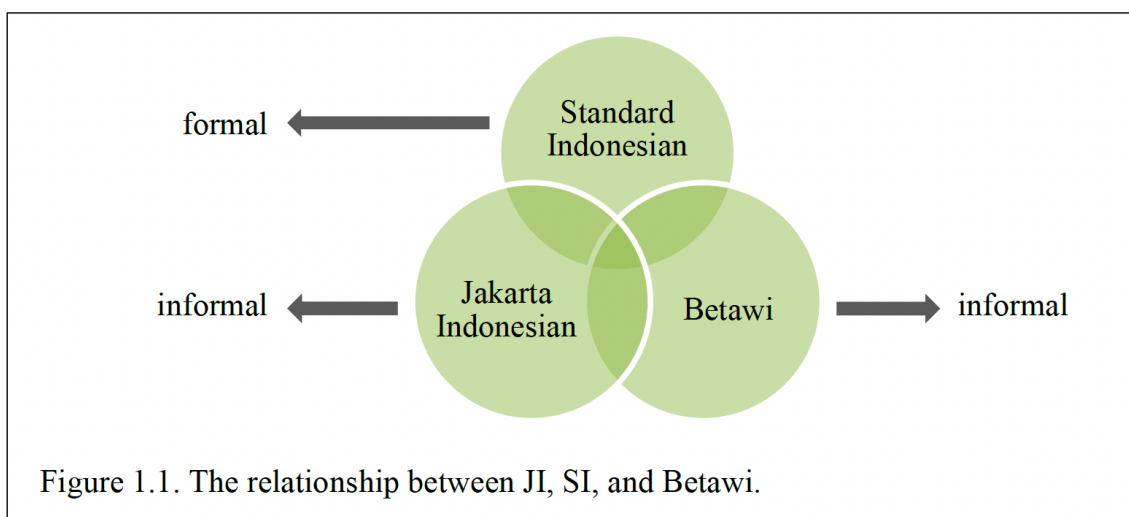
One of these colloquial varieties is Jakarta Indonesian (JI), the variety spoken in the nation's densely populated, urban capital of Jakarta. JI is an increasingly important variety of spoken Indonesian (Ewing 2005, Sneddon 2006, Kurniawan 2018). JI is considered a contact variety between Betawi, the local variety of Malay historically spoken in the area, and Standard Indonesian. Ikranagara (1975) observes:

As the national language takes over more informal and casual functions it is particularly Betawi which serves as a source for the developing casual lect. (p. 6)

As well as the borrowing of vocabulary, the influence of Betawi on Indonesian may be phonological or syntactic. (p. 9)

There is probably a continuum situation between Betawi and Bahasa Indonesia in Jakarta. (p. 11)

JI/Betawi and SI are in a diglossic relationship with JI or Betawi serving as an informal, largely spoken variety used together with SI as a more formal variety (Ikranagara 1975, Sneddon 2003b, Kurniawan 2018). This is schematized by Kurniawan (2018:3, Figure 1.1) repeated here as Figure 1.



**Figure 1. Kurniawan (2018:3), Figure 1.1**

As the native language of an increasingly large population of speakers in and around Jakarta, JI is not only used as a spoken variety in informal settings, but is also spread through print, audio-visual and social media. Jakarta Indonesian has exerted influence on other urban varieties through its use by “students and educated people” (Poedjosoedarmo 1982, Manns, Cole & Goebel 2016, Ewing 2019) moving between Jakarta and other cities, and dissemination by various forms of media. Not surprisingly, there is a lot of variation observed, but little prior study of the social and linguistic factors that condition that variation (but see Sneddon 2006; Kurniawan 2018).

## 2. Patterning of three phonological variables – Kurniawan (2018)

In this section we briefly review the groundbreaking work done by Ferdinan Kurniawan in his 2018 Cornell Ph.D. dissertation *Phonological Variation in Jakarta Indonesian: An Emerging Variety of Indonesian*.<sup>2</sup> We review the methodology, which we follow in our own work (section 2.1) and the findings for the three phonological variables investigated (sections 2.2-2.4). We then turn to conclusions and discussion (section 2.5), setting the stage for our follow up studies (section 2.6) reported in sections 3 & 4.

Kurniawan (2018: 182)

. . . investigates the development of Jakarta Indonesian (JI) using corpora based on three generations of naturalistic speech data to study variation in the realization of three (morpho-)phonological variables. . . . [I]t demonstrates the importance of naturalistic speech corpora in examining the actual patterns of language use focusing on colloquial speech, which we know to be the locus of language change. By studying naturalistic colloquial speech, it contributes to our understanding of linguistic variation, contact, and change in progress.

<sup>2</sup> In this section, we provide a summary of Kurniawan’s (2018) findings as presented in chapter 5, reproducing his figures and tables.

## 2.1 Methodology

Kurniawan (2018) studies three generations of speakers in an apparent time study, using the following corpora as summarized in (6):

(6) Three generations of JI speakers from corpus data:

- a. **Generation 1:** born 1945-1960 – Wallace (1976) available at <https://ecommons.cornell.edu/handle/1813/45606> (accessed 10/11/2022)
- b. **Generation 2:** born 1960-1980 – The Betawi-Jakarta Indonesian corpus Gil & Tadmor (2014) – adult speakers, available at [https://archive.mpi.nl/tla/islandora/object/tla%3A1839\\_00\\_0000\\_0000\\_0022\\_5AC9\\_0](https://archive.mpi.nl/tla/islandora/object/tla%3A1839_00_0000_0000_0022_5AC9_0) (accessed 10/11/2022)
- c. **Generation 3:** born 1990s – The MPI-EVA Jakarta Child Language Database Gil & Tadmor (2007) – pre-adolescents available at <https://childes.talkbank.org/access/EastAsian/Indonesian/Jakarta.html> (accessed 10/11/2022)

For the first generation, Kurniawan uses a corpus consisting of transcribed conversations from recordings collected in the mid-1970's under the direction of Stephen Wallace. The conversations were collected as part of Wallace's (1976) dissertation on phonological variation in JI (which he called Modern Jakarta Malay). A total of 35 hours of recordings involving over 200 adult speakers were collected by fifteen research assistants. These assistants recorded conversations among friends, relatives, and neighbors (Wallace was not present for the recordings). The recordings were transcribed in the 1970's after they were collected, and the transcriptions are now publicly available from Cornell University's digital repository through eCommons. The transcripts have been converted to machine-readable text and are searchable by character, word, and string of words.

For the second generation, Kurniawan uses the Betawi-Jakarta Indonesian corpus (the BJI corpus), and this was used in our follow up studies as well. The BJI corpus contains audio recordings of conversations collected between 2004–2012 under the auspices of the Max Planck Institute for Evolutionary Anthropology Jakarta field station (Gil & Tadmor 2014). The recordings were done in informal settings in Jakarta, with 143 speakers from a range of socio-economic and educational backgrounds with a total of 28 hours of recorded speech (a total of 75,079 utterances). The conversations in the corpus were segmented into utterances and transcribed by Indonesian linguists using ELAN. All utterances are translated, further segmented into morphemes, and glossed by morpheme. Thus, the corpus can be searched by speaker, conversation code, Indonesian word, morpheme, or morpheme gloss. For his study, Kurniawan analyzed the speech of 20 speakers who identified as JI speakers (using the same speakers for all three variables, see Kurniawan 2018:27, Table 1.3). The third generation, not the focus of our work here, uses the preadolescent data from the MPI-EVA Jakarta Child Language Database (see references cited above).

In each case, the target items are searched for in the relevant corpora and each token is analyzed based on the variant realized by speaker, age, gender, and educational level.

## 2.2 [-a] ~ [-e] variation

The first variable studied by Kurniawan (2018 chapter 2) is the alternation observed in JI word-final [-a] and [-e] corresponding to SI [-a]. The variants with final [-e] originated in Betawi in the inner city of Jakarta around the beginning of 19th century, coexisting with



the [-a] variant, the latter now interpreted as resulting from contact with SI. Kurniawan focuses on alternation in function words as these have the highest token frequency in the corpus and show the most variation.

(7) [-a] ~ [-e] variation examples

[ija] ~ [ije] 'yes'  
 [gua] ~ [gue] '1SG'  
 [apa] ~ [ape] 'what'

Kurniawan's (2018:184) findings are summarized in his Figure 5.1, reproduced here as Figure 2.

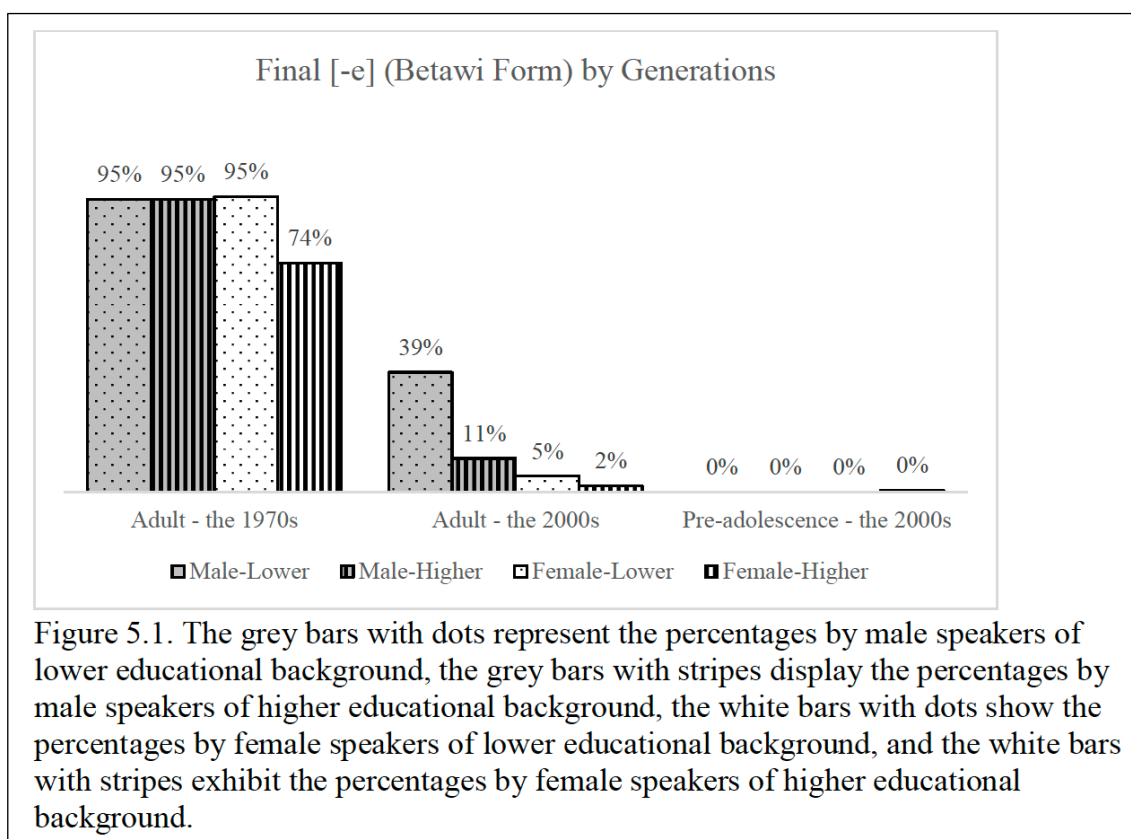


Figure 2. Kurniawan (2018:184), Figure 5.1

Kurniawan finds high rates of [-e] for all adult speakers in 1970s across the board (regardless of gender and education), however female speakers with high education exhibit lower rates than others. There are lower rates for all adults in the 2000s, though with male speakers with low education exhibiting the highest rate at this stage and with essentially no [-e] for preadolescent speakers in the 2000s. We clearly see variation conditioned by generation with highest rates of [-e] among adults in 1970s and lowest rates among pre-adolescents in 2000, reflecting an abrupt change from Betawi [-e] to SI [-a]. This is a classic example of change over time. Within this overall trend, there is some effect of gender and educational level, with evidence that women of higher educational level at the earliest stage are leading the change and with males of lower educational level in the mid stage lagging behind other groups.

**2.3 [-∅] ~ [-h] ~[-ʔ] variation**

In his second case study, Kurniawan (2018, chapter 3) examines the observed alternation among word-final [-∅], [-h], and [-ʔ] in function words in JI. Results presented are from the BJI corpus, since analyzing this pattern based on the transcriptions available in Wallace (1976) without audio would be unreliable.

- (8) [-∅] ~ [-h] ~[-ʔ] variation examples  
*iya* 'yes' [ija] ~ [ijah] ~ [ijaʔ]  
*lagi* 'more' [lagi] ~ [lagih] ~ [lagiʔ]

The [-ʔ] and [-h] forms are taken to be Betawi and Sundanese influenced forms respectively, with the [-∅] forms being evidence of the SI form. The envelope of variation is phonologically conditioned as the variation occurs phrase finally, but not phrase medially. The results of this variation as observed in phrase final position are summarized in Kurniawan’s (2018:185) Figure 5.2, reproduced here as Figure 3.

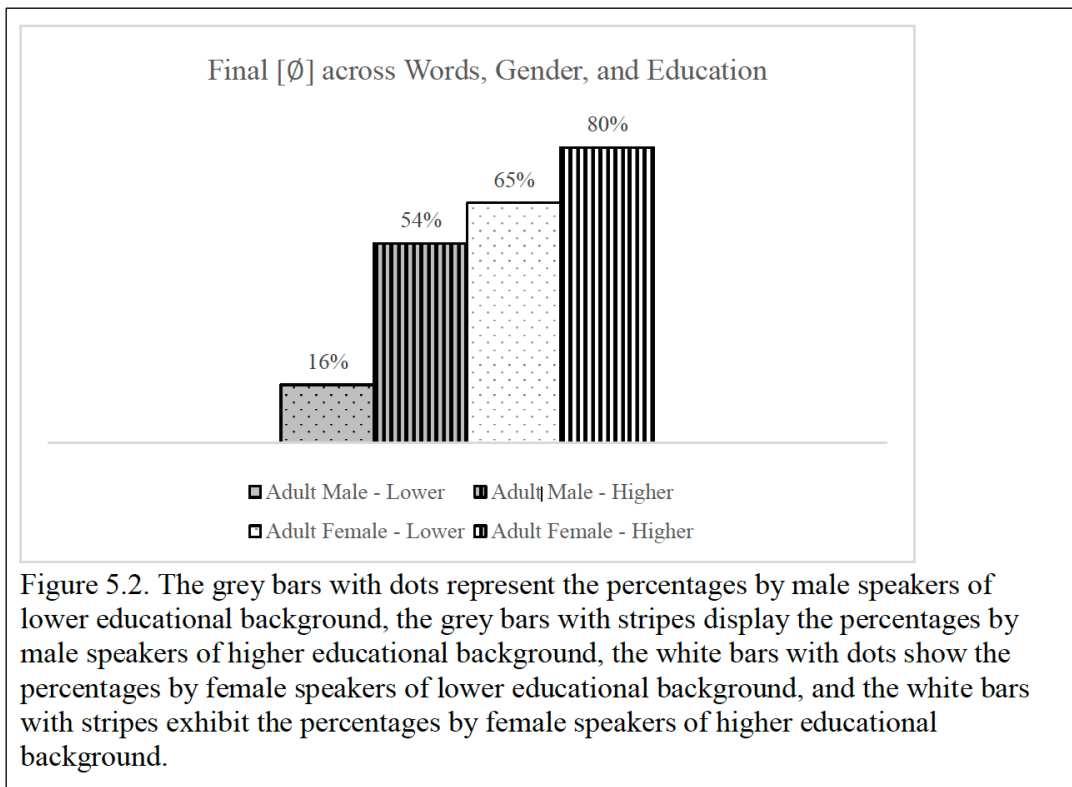


Figure 5.2. The grey bars with dots represent the percentages by male speakers of lower educational background, the grey bars with stripes display the percentages by male speakers of higher educational background, the white bars with dots show the percentages by female speakers of lower educational background, and the white bars with stripes exhibit the percentages by female speakers of higher educational background.

**Figure 3. Kurniawan (2018:185), Figure 5.2**

Kurniawan finds that the Betawi influenced forms are still present, but mixed with SI variant [-∅]. The rates of [-∅] are influenced by both gender and educational background with stepwise results: male lower education < male higher education < female low education < female higher education. Thus there is a clear shift toward the SI forms, but with the Betawi and Sundanese influenced forms still present and conditioned by both gender and educational level.

**2.4 ηə- ~ N- ~ bare verb ~ məN-**

In his third case study, Kurniawan (2018, chapter 4) examines the observed alternation in the realization of active verbs in JI. In SI, the prefix /məN-/ is expected, with the nasal

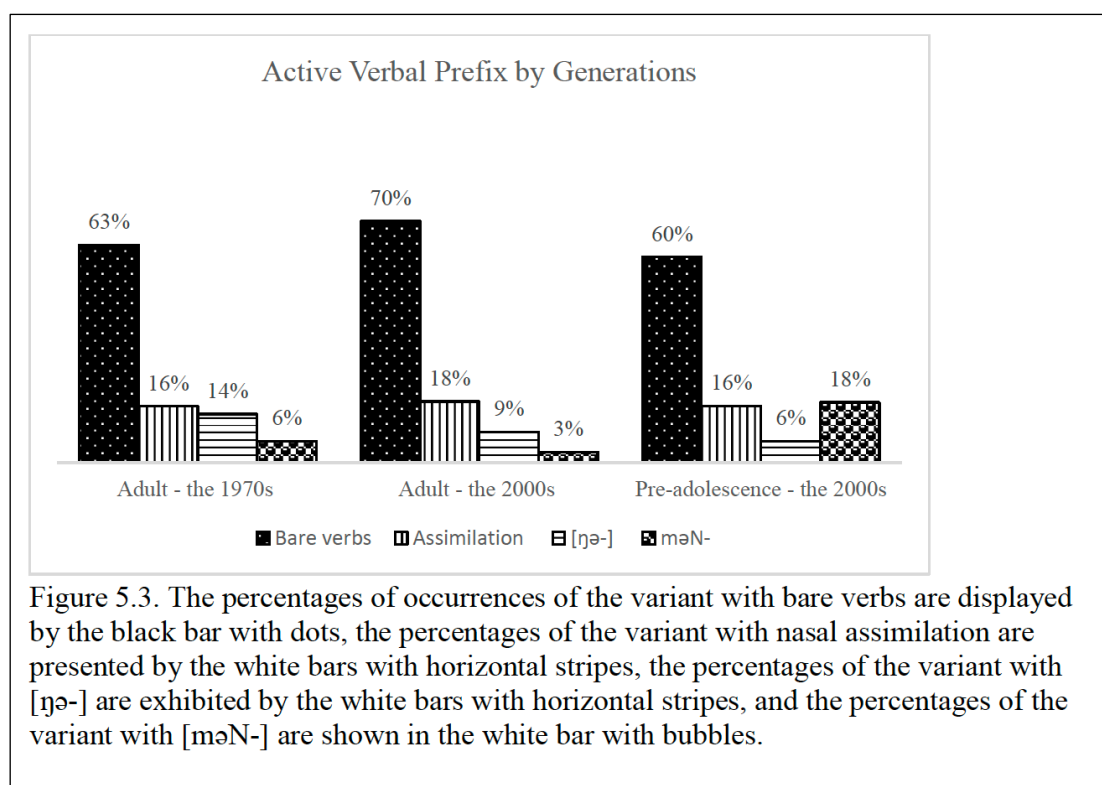
surfacing as [ŋ] before vowel initial roots, assimilating in place of articulation to verb root initial obstruents (and assimilation and deletion in the case of voiceless obstruents) and deleting before sonorants (see Lapoliwa (1981) for a comprehensive treatment). In JI, the variants ~ [ŋə-] ~ [N-] are both expected, phonologically conditioned by the root initial sound, expected to surface as [ŋə] before sonorants, as [ŋ] before vowels and with assimilation and deletion with voiceless obstruents. In the case of voiced obstruent initial roots, variation is observed with both forms occurring. In addition to these two forms, Kurniawan found rather unexpectedly that it is also common for a bare verb form (just the root, no prefix) to be used as well as the SI [mən-] form. We look at the overall distribution before focusing more specifically on the [ŋə-] ~ [N-] alternation, observed for verb stems starting with voiced obstruents.

(9) ŋə- ~ N- ~ mən- ~ bare verb examples

Jl variation for active form of /bəli/ 'to buy':

∅-bəli ~ ŋə-bəli ~ m-bəli ~ məm-bəli

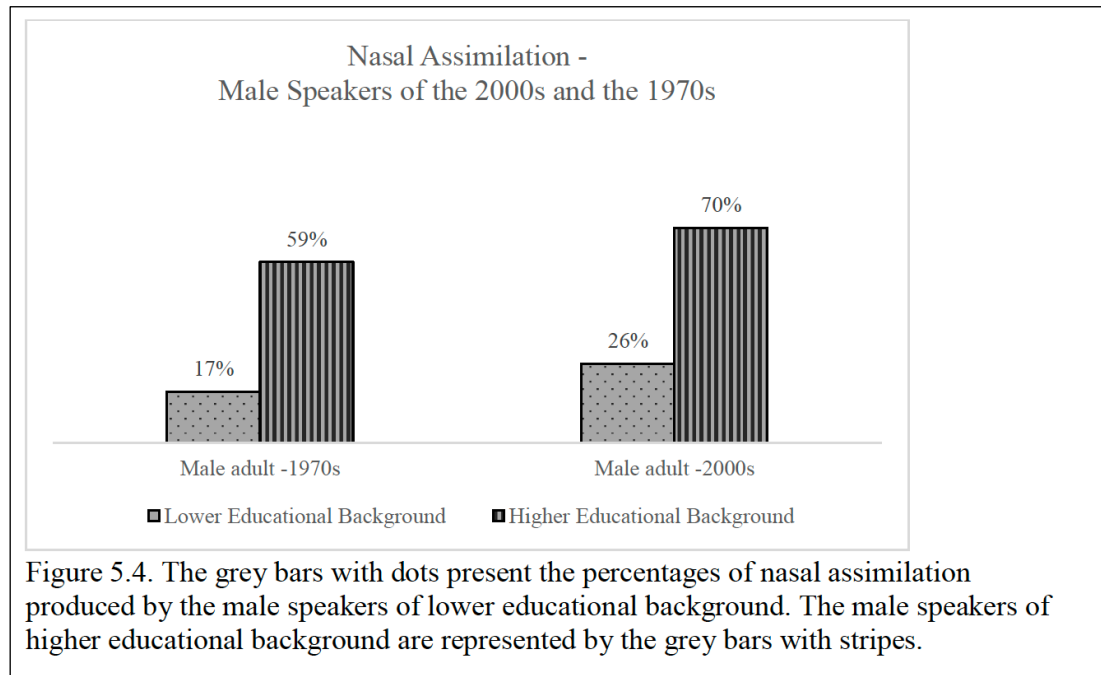
The results of this variation are summarized in Kurniawan (2018:186) Figure 5.3, reproduced here as Figure 4.



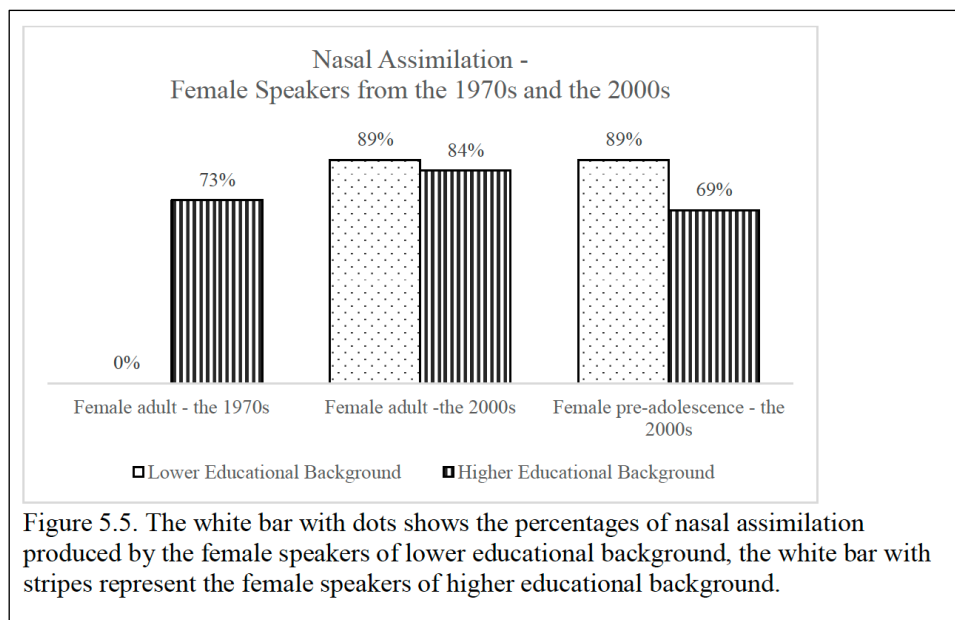
**Figure 4. Kurniawan (2018:186), Figure 5.3**

The very high use of bare verb forms in contexts that clearly involve the active verb form is unexpected linguistically and warrants further investigation, but does not seem to be conditioned by social factors as these forms are used equally commonly across age, gender, and education level. There is a minimal use of [mən-] forms among adults both in the 1970s and 2000s which Kurniawan interprets as code switching. He further notes an increase among the preadolescents suggesting this might indicate an integration of this form into the colloquial register.

Turning to a closer look at  $[\eta\text{ə-}] \sim [N-]$ , seen with voiced initial roots (Kurniawan 2018:145, example (11)), the results of this variation are summarized by Kurniawan 2018:188–189, Figure 5.4 & 4.5, reproduced here as Figures 5 & 6.



**Figure 5. Kurniawan 2018:187, Figure 5.4**



**Figure 6. Kurniawan 2018:188, Figure 5.5**

Overall, there is a greater difference between education levels for men than for women with less educated male speakers making much greater use of the  $[\eta\text{ə-}]$  form than the nasal assimilated form, as shown in Figure 5. As shown in Figure 6, for female speakers, this difference is only seen among the speakers from the 1970s. Thus there is a demonstrated change in patterns of variation seen between the 1970s to 2000s. Insufficient data were

available for the preadolescent males, but the pattern for females is basically the same as for the 2000s adult females.

## 2.5 Conclusions and discussion

Kurniawan (2018:183) concludes:

The findings from the three variables . . . support the conclusion that JI is in fact an admixture of Betawi and SI, with a strong influence of Javanese, Sundanese, and Bangka Malay. The relationship between these varieties is identified in the patterns of variation that show a general trend toward increased use of the SI and Javanese variants.

Kurniawan (2018:189) provides a comparison of the three variables in Table 5.1, reproduced here as Table 2, comparing the roles of age (as grouped by generation), linguistic conditioning, gender and education, suggesting three different specific patterns of shift, while all showing the same overall trend. He observes “Notably, we can see that the social factors that condition the variation are not the same for all variables” (p. 189).

**Table 2: Kurniawan (2018:189), Table 5.1**

	Variables								
	[-a] ~ [-e]			Ø ~ [-h] ~ [-ʔ]			Assimilation ~ [ŋə-]		
Generations	A-70s	A-2000s	PreA-2000s	A-70s	A-2000s	Pre-A-2000s	A-70s	A-2000s	Pre-A-2000s
Linguistic Conditioning	X	X	X	N/A	Phrase Boundary		X	X	X
Gender	X	X	X	N/A	√	X	√	√	X (?)
Education	X (?)	X (?)	X	N/A	√	X	√	√ (M)	X
Major Points	Abrupt Change: Betawi to SI			Moving towards SI, but highly determined by gender and educational background			<ul style="list-style-type: none"> <li>• Moving towards Javanese form</li> <li>• Style-shifting between standard and colloquial forms</li> </ul>		

The fact that the overall trends are the same but each variable tells a somewhat different story is a very interesting result which warrants fuller investigation. Kurniawan (2018:193) concludes:

Further investigation on other linguistic variables in JI is needed to see if the results we have in this study are also applied to other variables. This would allow us to see if the degree to which different variables are used to indicate similar socio-indexical effects. Do other variables work the same way?

Kurniawan (2018:194) suggests that “the influence of (H) on (L) is less likely to occur at

the morpho-phonological level than at the phonological level. To answer this, we need to further study the morphological and syntactic variables in the corpora.”

## 2.6 Next steps

In our current project, we take up this invitation to investigate additional variables for which data are available in the BJI corpus and to the degree possible looking at the same speakers investigated by Kurniawan. We hope to contribute to the question of how the patterning of different variables is related; and to understand how multiple patterns of variation in the same language fit together. How do different conditioning factors contribute to specific patterns of variation? How do linguistic and discourse factors interact with social context and identity of speakers or interlocutors? Does the variation tell us something about the speakers themselves or about the context they are speaking in?

In sections 3 and 4, we turn to the investigation of three additional linguistic variables in the Betawi-Jakarta Indonesian corpus. These are 1) alternation between a word-initial consonant and  $\emptyset$  (zero), based on Cohn & Vogel (2019), and 2) lexical variation involving the first person singular pronoun, based on Abtahian, Cohn, Djengar & Vogel (2021). (Note that the latter study also investigated variation within the BJI corpus for both JI and Betawi speakers and within the Wallace corpus; however, we focus here on the findings only in the BJI corpus for JI speakers.) For each variable, we examine patterns of inter- and intra-speaker variation and consider to what extent the observed variation is conditioned by various linguistic factors, social factors, and/or by speech style. Crucially, we study these patterns in the same set of speakers for all three variables, allowing for a meaningful comparison both between different types of variables (lexical vs. phonological) and between different variables of the same type (two superficially similar phonological alternations).

## 3. Word-initial consonant ~ $\emptyset$ alternations

The two word-initial alternations we examine are between [h] ~  $\emptyset$  and [s] ~  $\emptyset$ . While we might expect these two to be similar given that there is an alternation between the presence or absence of a word initial consonant in each case, prior descriptions suggest that they pattern differently in Jakarta Indonesian. In what follows, we start by reviewing the prior literature on the two alternations and then present the results of our investigation in the Betawi-Jakarta Indonesian corpus.<sup>3</sup>

### 3.1. Prior descriptions of word-initial alternations between [h] ~ $\emptyset$ and [s] ~ $\emptyset$

The alternation between word-initial [h] and  $\emptyset$  is illustrated in (11a–c) below. This alternation is not said to be associated with particular lexical items but is instead thought of as an optional phonological rule of deletion (e.g., /h/  $\rightarrow$   $\emptyset$  / #\_\_\_) historically related to loss of initial /h/ in related varieties of Malay (Ikranagara 1980, Ewing 2005, Sneddon 2006).

(10) Word-initial [h] ~  $\emptyset$  alternation

---

<sup>3</sup> In this section, we provide a summary of Cohn and Vogel’s (2019) findings, reproducing their figures and tables.

- a. hari ~ ari 'day'  
 b. habis ~ abis 'finished'  
 c. hijau ~ ijo 'green'

In contrast, the alternation between word-initial [s] and  $\emptyset$ , illustrated in (11a–c), has been described as being associated with certain lexical items, in particular, high frequency items and grammaticalized versions of lexical forms (Ewing 2005, Sneddon 2006). For example, while (11b) shows an alternation between *sampe* and *ampe*, the version without the [s] is only available for the grammatical form meaning 'until,' not for the lexical form meaning 'arrive.'

(11) Word-initial [s] ~  $\emptyset$  alternation

- a. saja ~ aja 'just'  
 b. sampe ~ ampe 'until' (cf. sampe ~ \*ampe 'arrive')  
 c. suda(h) ~ uda 'perfective'

Additionally, the word-initial [s] ~  $\emptyset$  alternation has been found to be associated specifically with casual speech, as a marker of informal register. Sneddon (2006:19, Table 1a & 1b), repeated here as Table 3, illustrates this effect. As can be seen, in the data from conversational and interview settings, Sneddon found that nearly all tokens of both *aja/saja* and *udah/sudah* occur without the initial [s]. The [s]-initial variants are almost entirely restricted to meeting settings.

**Table 3. Frequency of [s] vs.  $\emptyset$  -initial forms in two lexical items by speech setting (reproduced from Sneddon 2006:19 Tables 1a & 1b)**

<b>Table 1a:</b> Frequency of <i>aja/saja</i> variants				
	<i>aja</i>	<i>saja</i>	total	% <i>aja</i>
conversations	505	6	511	98.8
interviews	401	26	427	93.9
	45	37	82	54.9

<b>Table 1b:</b> Frequency of <i>udah/sudah</i> variants				
	<i>udah</i>	<i>sudah</i>	total	% <i>udah</i>
conversations	737	31	768	96.0
interviews	803	85	888	90.4
	33	74	107	30.8

Given the divergent descriptions of the two word-initial consonant ~  $\emptyset$  alternations, we predict that they would be conditioned by different factors within the Betawi-Jakarta Indonesian corpus. Specifically, we predicted that the [s] ~  $\emptyset$  alternation would exhibit a substantial effect of lexical frequency. It should be noted that since the Betawi-Jakarta Indonesian corpus only contains casual speech, we were not able to directly investigate

the register effect observed in Sneddon (2006), but rather we investigated additional variation within the casual speech style. We also predicted that the [h] ~ Ø alternation (but not necessarily the [s] ~ Ø alternation) would be conditioned by phonological factors, in particular the environment to the left of the word in question (i.e., the right edge of the preceding word).

### 3.2. [h] ~ Ø and [s] ~ Ø in the Betawi-Jakarta Indonesian corpus

The data examined in our investigation of the two word-initial consonant ~ Ø alternations consisted of all /h/-lexemes in the corpus exhibiting initial alternations with at least five tokens in the corpus and all the /s/-lexemes in the corpus exhibiting initial alternations. The data were then filtered down to include only the tokens of these lexemes that were produced by the specific subset of speakers included in Kurniawan’s study (see Kurniawan 2018:27–28, Table 1.3).

Table 4 summarizes the results for each lexeme-type. Note that in the rest of this section, we refer to the forms without the initial consonant as “vowel-initial” or “V-initial” rather than “Ø-initial” in order to reflect the possibility that these forms exist independently in speakers’ lexicon rather than resulting only from active phonological deletion. The forms with the initial consonant are referred to as “consonant-initial” or “C-initial.”

**Table 4. Rates of C-initial and V-initial forms by lexeme type**

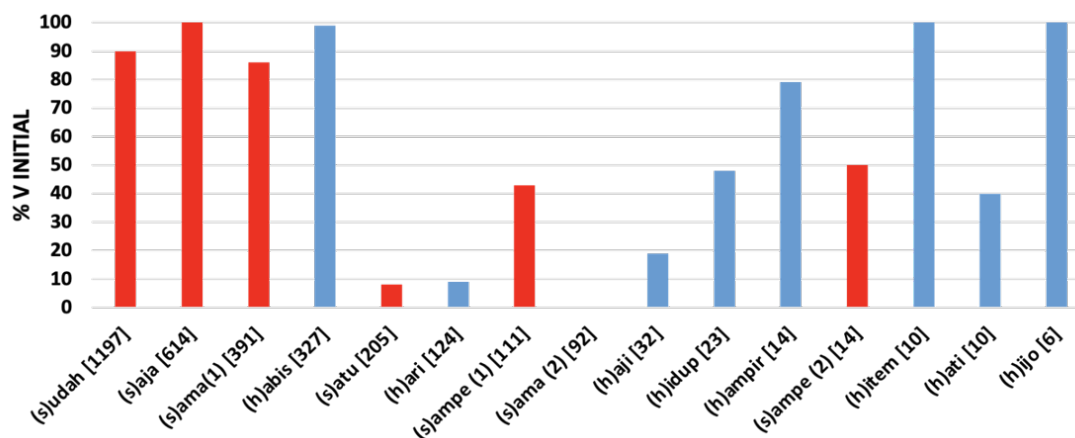
	<b>C-initial</b>	<b>V-initial</b>	<b>Total</b>
<b>/h/-lexemes</b>	46%	54%	362
<b>/s/-lexemes</b>	12%	88%	2532

As can be seen in this table, there is a roughly even split between consonant-initial and vowel-initial tokens for the /h/ lexemes (46% C-initial; 54% V-initial). In contrast, the tokens of /s/-lexemes are almost all vowel-initial (12% C-initial; 88% V-initial). Since, as noted above, the Betawi-Jakarta Indonesian corpus contains only casual speech, the very high rates of V-initial forms for the /s/-lexemes in particular is precisely what we should expect following Sneddon (2006). The rest of this subsection looks more closely at the variation within both lexeme types, starting by breaking the data down according to the individual lexical item, then turning to phonological conditioning and social conditioning.

#### 3.2.1 Variation across lexemes

Figure 7 presents a closer look at variation across individual lexical items and the effect of token frequency on this variation. Each bar in this figure represents an individual lexical item in our data (/s/-lexemes with grammaticalized and lexical forms are split into two separate bars), and they are ordered on the x-axis from highest token frequency to lowest token frequency.





**Figure 7. Results by [lexeme frequency] (highest → lowest token frequency):**  
 red = /s/-lexemes, blue = /h/-lexemes; *sampe(1)* is ‘arrive’,  
*sampe(2)* is ‘until’; *sama(1)* is ‘with’, *sama(2)* is ‘same’

Recall that prior descriptions of the [s] ~ ∅ alternation suggest that it is associated with high frequency and grammaticalized forms, whereas the [h] ~ ∅ alternation is not described as being restricted to particular lexical items. For /s/-lexemes, this figure is largely consistent with the previous descriptions. We see that the rates of V-initial forms for higher frequency /s/-lexemes (the red bars closer to the left side of the x-axis) are above 80%, whereas the rates for lower frequency /s/ items are at or below 50%. In terms of the effect of grammaticalization, the pattern for *sama/ama* is consistent with the previous descriptions as well. That is, the grammaticalized version, meaning ‘with,’ is almost always V-initial (see *(s)ama(1)* in the figure), whereas the lexical version, meaning ‘same,’ is always C-initial (see *(s)ama(2)* in the figure). The pattern for *sampe/ampe*, however, does not show an effect of grammaticalization. That is, both grammaticalized and lexical versions exhibit roughly comparable percentages (40–50%) of V-initial tokens (compare *(s)ampe(1)* and *(s)ampe(2)* in the figure).

For /h/-lexemes, as predicted, we do not find an effect of frequency, since both high and low frequency items can occur in their V-initial forms more than 90% of the time (e.g., high frequency *(h)abis* and low frequency *(h)item* and *(h)ijo*). We do, however, still see substantial variation across lexical items for the /h/-lexeme category, although it is not conditioned by frequency. That is, while some lexemes occur almost entirely in their V-initial form, as noted as noted above, others exhibit very low rates of %V-initial (e.g., *(h)ari* and *(h)aji*).

### 3.2.2 Phonological conditioning

Next, we investigated the role of phonological conditioning, testing the prediction that variation in the /h/-lexeme category, but not necessarily in the /s/-lexeme category, is conditioned by the phonological environment to the left of the word in question. In particular, we predicted that if the alternation is phonologically conditioned, V-initial forms would be more likely after consonant-final words and C-initial forms would be more likely after vowel-final words. (In this scenario, either deletion would occur to avoid consonant clusters, or it would be blocked to avoid vowel hiatus.) This prediction was not borne out, however. Tables 5 and 6 below illustrate our results for two lexemes we examined—one /h/-lexeme, *hidup* ‘live’, and one /s/-lexeme, *sampe* ‘until’.

**Table 5. Rates of C-initial and V-initial variants for (*h*)*idup* depending on phonological environment. Gray shaded cells predicted to exhibit higher rates; white cells predicted to exhibit lower rates)**

<b>hidup</b>	<b>C # _</b>	<b>V # _</b>
C-initial	55%	45%
V-initial	56%	44%

**Table 6. Rates of C-initial and V-initial variants for (*s*)*ampe*(2) depending on phonological environment. Gray shaded cells predicted to exhibit higher rates; white cells predicted to exhibit lower rates)**

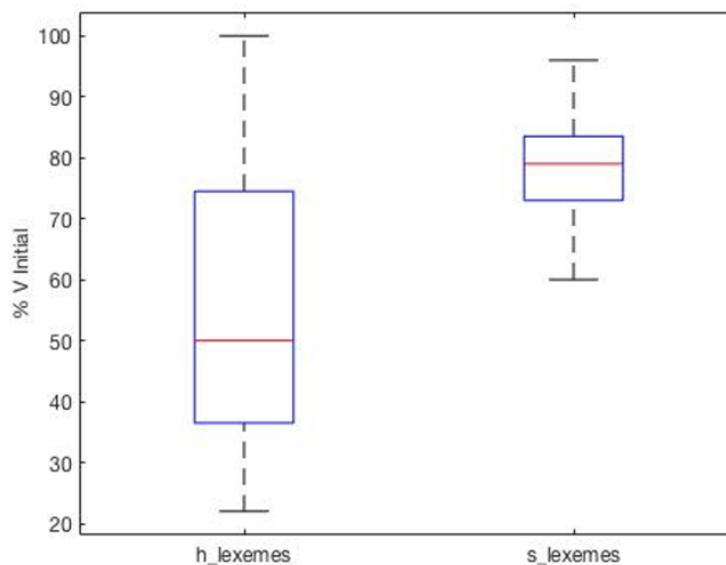
<b>sampe (2)</b>	<b>C # _</b>	<b>V # _</b>
C-initial	71%	29%
V-initial	53%	47%

As these tables show, for (*h*)*idup* and (*s*)*ampe*, the environments in which we predicted higher rates of the V-initial variant exhibit roughly even rates of V-initial and C-initial variants, or even higher rates of the C-initial variant. In contrast, the environments in which we predicted lower rates of the V-initial variant exhibit either comparable or higher rates to both the C-initial variant in the same environment and the V-initial variant in the other environment. A linear regression analysis that treated % *V-initial* as the response variable and *phonological environment* and *lexeme type* as main and interacting fixed effects verified that phonological environment was not a significant predictor of variation in our data ( $p > 0.05$ ).

This is not surprising for  $s \sim \emptyset$  described as lexicalized and as shown above having token frequency effects. However it is surprising for  $h \sim \emptyset$  which has been described as an optional phonological rule and therefore we would expect to be sensitive to phonological conditioning.

### 3.2.3 Intraspeaker variation, interspeaker variation, and social conditioning factors

Next, we considered to what extent the overall variation observed between C-initial and V-initial forms results from interspeaker differences and/or intraspeaker variation and whether the interspeaker variation is conditioned by social characteristics of the speakers. Figure 8 shows rates of V-initial forms produced by each speaker for /h/-lexemes and /s/-lexemes separately. (% V-initial reflects the percentage of V-initial forms a speaker produced out of all their tokens of the relevant lexeme category.)



**Figure 8. Rates of V-initial forms of /h/-lexemes and /s/-lexemes for each speaker**

This figure shows a particularly high degree of interspeaker variation for the /h/-lexemes, with some speakers producing around 20% V-initial /h/-lexemes and others producing 100% V-initial /h/-lexemes. For the /s/-lexemes, no speaker produces less than around 60% V-initial. This is consistent with the overall higher rates of V-initial variants found for the /s/-lexemes relative to the /h/-lexemes in the corpus (see Table 4), and with prior descriptions of the V-initial variant for /s/-lexemes being associated with casual speech. Nevertheless, there is still substantial interspeaker variation even for /s/-lexemes. While all speakers produce majority V-initial forms for the /s/-lexemes, their rates of V-initial forms range from ~60 to ~100%, indicating that different speakers exhibit distinct patterns. This figure also reveals substantial intraspeaker variation, especially for /h/-lexemes. While some speakers produce 90-100% V-initial forms, many other speakers produce midrange percentages of V-initial forms (e.g., 40–60%). These midrange percentages indicate that a single speaker produces both C- and V-initial variants at comparable rates.

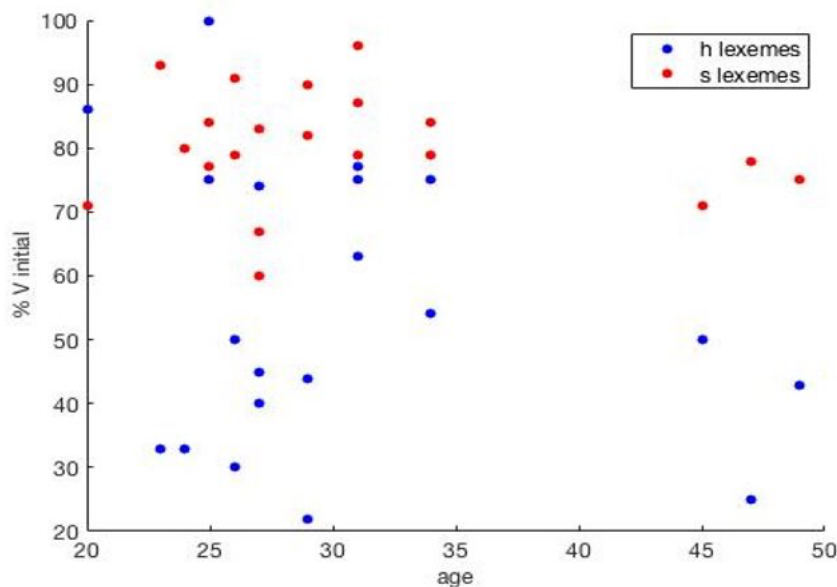
Looking more closely at the variation seen across speakers, we predicted that it was conditioned at least to some degree by speakers' social characteristics. Specifically, we predicted that V-initial forms would be more likely for males, for speakers with lower education levels, and for older speakers, because these groups should have less influence from Standard Indonesian, in which only the C-initial forms are possible. Table 7 shows the rates of V-initial forms broken down by gender, education level, and lexeme type.

**Table 7. Rates of V-initial /h/ and /s/-lexemes for female vs. male speakers and for speakers with lower vs. higher levels of educational attainment**

	% V /h/-lexemes	% V /s/-lexemes
<b>Female</b>	56%	77%
<b>Male</b>	54%	80%
<b>Lower Ed.</b>	57%	86%
<b>Higher Ed.</b>	55%	88%

Contrary to our predictions, this table shows that within each lexeme type, there are comparable rates of V-initial forms between genders and between lower and higher levels of educational attainment. *t*-tests verified these findings: no significant difference was found between genders or between education levels.

Figure 9 shows the % V-initial for each speaker and each lexeme type with speakers ordered from youngest to oldest on the x-axis.



**Figure 9. %V-initial for each speaker and lexeme type with speakers ordered from youngest to oldest. Red dots are %V-initial for /s/-lexemes; blue dots are %V-initial for /h/-lexemes**

As this figure shows, there does not appear to be a correlation between age and %V-initial. Indeed, the correlation coefficient was not significant for either lexeme type ( $r = -0.25$  for /h/-lexemes and  $-0.15$  for /s/-lexemes;  $p > 0.05$  for both). Thus, while there is substantial interspeaker variation in the rates of %V-initial for both /h/-lexemes and /s/-lexemes, as seen in Figure 8, this variation is not conditioned by any of the social factors that we examined. Figure 9 also shows, consistently with Figure 8, that many speakers produce comparable rates of V-initial and C-initial variants, reflecting substantial intraspeaker variation. Thus while we do observe considerable inter- and intra-speaker variation, this does not appear to be socially conditioned.

#### 4. First person singular pronoun variation

The third variable we examined was lexical variation in first person singular pronouns (abbreviated 1SG in the rest of this section). Jakarta Indonesian includes a relatively large number of forms for self-address, all corresponding to “I / me” in English. These include *saya* and *aku* from Standard Indonesian, as well as borrowings from several other sources, and are said to have different social and pragmatic associations. We start by reviewing prior descriptions of major variants in the literature and then present the results of our

investigation of the variation in 1SG pronouns in the Betawi-Jakarta Indonesian corpus.<sup>4</sup>

#### 4.1. Prior descriptions of 1SG variants

Table 8 lists seven 1SG variants along with the previous descriptions of their historical sources and social and pragmatic associations.

**Table 8. Variants of 1SG pronoun under examination, based on Djenar et al. 2018, Englebretson 2007, Ewing 2005, Ewing 2019, Manns 2014, Sneddon 2006**

Pronoun	Previous descriptions
saya	Sanskrit borrowing; public identity; used in interactions between those who are not social intimates
gue	Hokkien borrowing via Betawi; associated with Jakartan youth identity. In Bandung, associated with outspokenness, exaggerated speech, and bravado
gua	Hokkien borrowing, common among Jakarta speakers, less linked with youth identity than <i>gue</i> ; commonly used by ethnic Chinese speakers in other cities
aku	From Malay; used in interactions between social intimates; indexes personal identity, more relaxed and intimate self
-ku	Clitic form of <i>aku</i>
kita	1PL inclusive in Standard Indonesian, inclusive and exclusive in colloquial Indonesian including Jakarta Indonesian; also used to denote first person singular
kite	Betawi Malay variant of <i>kita</i> , mainly used as 1SG

As discussed by Abtahian, Cohn, Djenar & Vogel (2021), the fact that such a wide variety of pronominal forms is used is of considerable interest. For our purposes here, it is the wide range of both inter- and intra-speaker variation observed that is our focus.

#### 4.2. 1SG variation in the Betawi-Jakarta Indonesian corpus

In order to investigate 1SG variation, we extracted all overt tokens of a 1SG pronoun produced by the 40 speakers under investigation (20 JI and 20 Betawi, including a subset of the Betawi speakers studied by Kurniawan 2015). We found that the vast majority of these tokens were *saya*, *gua*, and *gue* (3,685 out of 3,795 total overt 1SG tokens) and therefore focused our analysis on these three variants, since the large amount of data was ideal for quantitative analysis. We also ultimately treated *gua* and *gue* as a single category, given the similarity between the two forms seen in Table 8 and the fact that they are sometimes understood as part of the variable process of /a/ → [e] raising (e.g., Kurniawan 2018). Table 9 presents the overall rates of *saya*, *gua/gue*, and other variants in the data.

---

<sup>4</sup> In this section, we provide a summary of Abtahian, Cohn, Djenar & Vogel (2021) findings, reproducing their figures and tables.

**Table 9. Frequency of 1SG variants**

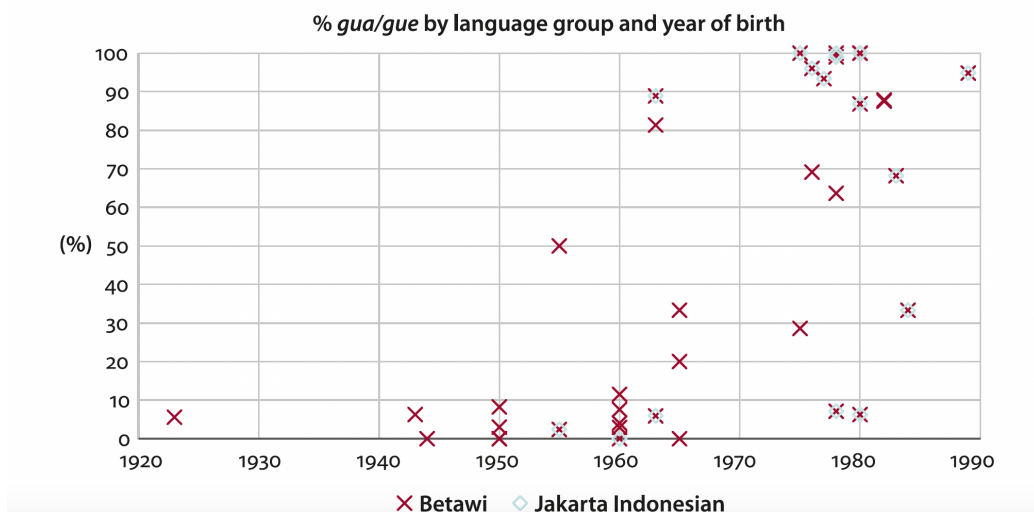
Variant	Number of tokens	% of total overt 1SG tokens
<i>saya</i>	2,160	56.9%
<i>gua/gue</i>	1,525	40.2%
<i>other</i>	110	2.9%

#### 4.2.1 Interspeaker variation and social conditioning factors

We first filtered down the data from the corpus to include only the 1SG tokens that were *saya* or *gua/gue* and calculated the % *gua/gue* for each of the 40 speakers. We found that this ranged from 0 to 100%, reflecting a large degree of interspeaker variation. In order to determine to what extent this variation across speakers was conditioned by their social characteristics, we conducted a mixed-effects logistic regression analysis with pronoun variant (*saya* vs. *gua/gue*) as the response variable, the social factors gender, birth year, and language background as fixed effects, and speaker as a random effect. Note that we also centered and scaled birth year by subtracting the mean and dividing by the standard deviation.

Given the previous descriptions of *gua* and *gue* as being associated with Jakarta and, in particular, Jakartan youth identity (e.g., Manns 2014; Djenar, Ewing & Manns 2018; Ewing 2019), and the strong association of Betawi Malay with Jakarta (e.g., Ikrangara 1975; Sneddon 2006; Kurniawan 2015) we predicted to find significant effects of birth year and language background. Specifically, the expectation was that younger speakers and those with a Betawi language background would use higher rates of *gua/gue*, whereas older speakers and those with a JI language background would use higher rates of *saya*. Our prediction for birth year was borne out ( $p < 0.001$ ) but not for language background ( $p > 0.05$ ). Gender was also not found to be significant ( $p > 0.05$ ).

Figure 10 shows the effect of birth year on the pronoun usage of individual speakers, where each of the 40 speakers is plotted separately. The x-axis corresponds to their birth year and the y-axis to their % *gua/gue*. As can be seen, as birth year increases (i.e., for younger speakers), their rate of % *gua/gue* also increases. This figure also indicates the language background of the speakers.



**Figure 10.** % *gua/gue* for each speaker ordered from earliest to latest birth year on the x-axis. Red Xs are Betawi speakers; blue diamonds are JI speakers.

Figure 10 also illustrates another important point about the variation in the rates of *gua/gue* among the 40 speakers. As noted earlier in this section, some speakers produce 0% *gua/gue* and others produce 100% *gua/gue*. These speakers are categorical users of either *saya* or *gua/gue*. At the same time, however, it can be seen in this figure that many speakers produce somewhere between ~30% and ~70% *gua/gue*. These speakers also produce between 30% and 70% *saya*, meaning that they exhibit intraspeaker variation.

#### 4.2.2 Intraspeaker variation

Given the large degree of intraspeaker variation between *gua/gue* and *saya*, we conducted a more in depth investigation of the variation within a single speaker in the corpus (male, JI language background, born in 1977), to determine whether his variation was conditioned by conversation-level factors. This speaker, who is coded as EXPOKK in the corpus, produced the most 1SG tokens out of all 40 speakers we examined and participated in 14 conversations, providing enough data for robust patterns to emerge, and allowing us to compare rates of *gua/gue* across a large number of conversations.

First, we found that the number of 1SG tokens produced by EXPOKK ranged from none to 179 across the 14 conversations. Second, looking just at the conversations in which he produced at least some 1SG tokens, his rates of both *gua/gue* and *saya* ranged from 0–100%, reflecting major variation across the conversations. At the same time, however, there were very few conversations in which EXPOKK produced midrange percentages of either pronoun variant. Rather, he typically produced a given variant either 0 to ~10% of the time or ~90 to 100% of the time within a single conversation. This suggests that whether EXPOKK used *gua/gue* or *saya* was largely set at the conversation-level.

We then investigated whether EXPOKK's choice of pronoun variant was conditioned by his interlocutors. However, we found that among the conversations in which he used high rates of *gua/gue*, some included conversations in which the other participants also used high rates of *gua/gue* as well as ones in which the other participants used high rates of *saya*. This indicates that he did not necessary select a pronoun variant to conform to the patterns produced by interlocutors.

Instead, we found that the most striking factor in EXPOKK's choice of pronoun variant

in a given conversation seems to be the total number of 1SG tokens he produced. Specifically, in conversations in which he produced large numbers of 1SG tokens, they were 90–100% *gua/gue*. In conversations in which he produced fewer than 10 1SG tokens, on the other hand, he used much higher rates of *saya*. Crucially, we also saw that in the conversations in which EXPOKK used 0–10 1SG tokens, the other participants often produced much larger numbers of 1SG tokens. We found that EXPOKK played a different role in the low and high % *gua/gue* conversations. He was a research assistant in the development of this corpus, and in cases in which he produced very few 1SG tokens, he was acting as an interviewer or facilitator of the conversation, sharing little about himself but asking other participants to talk about themselves. In contrast, when he produced larger numbers of 1SG tokens, he was acting as a regular conversation participant, talking about himself just as other participants talk about themselves. We therefore conclude that EXPOKK's variation is largely conditioned by his role in the conversation. When he is a regular conversation participant, he uses primarily *gua/gue*, and when he is an interviewer or facilitator, he uses primarily *saya*, consistent with the public association of *saya* mentioned in previous research noted above, compared to a more personal/interpersonal use of *gua/gue*.

### 4.3. Summary and comparison across the three variables

Sections 3 & 4 have presented the results of our investigation of three linguistic variables in the Betawi-Jakarta Indonesia corpus, including two that appear at first glance to be phonological (i.e., alternations between word-initial consonants and  $\emptyset$ ) and one lexical variable (i.e., the form of the 1SG pronoun). This investigation revealed substantial variation both across the set of speakers examined and within the speech of individual speakers, corresponding to inter- and intraspeaker variation respectively. At the same time, we found that the factors conditioning each level of variation were different for the different variables. Specifically, the interspeaker variation seen for the word-initial consonant  $\sim \emptyset$  alternations was not conditioned by any of the social factors we considered, including age, whereas age was a significant predictor of the interspeaker variation in the 1SG pronouns. In terms of intraspeaker variation, we found the role of a speaker in a given conversation to be an important factor for 1SG pronouns, whereas lexical frequency appears to be important for at least some of the variation in word-initial consonant  $\sim \emptyset$  alternations. Even within the two superficially similar patterns involving initial consonants, however, frequency only seemed to influence the [s]  $\sim \emptyset$  alternation—not the [h]  $\sim \emptyset$  alternation, consistent with the diverging prior descriptions of these two alternations and demonstrating the importance of studying and comparing across multiple patterns of variation within a single language community.

## 5. Conclusions and discussion

We now return to our questions above about how different variables might pattern in terms of both their linguistic and social conditioning. In Table 10 we compare Kurniawan (2018)'s result for the three variables he studied, with the three variables we add to the discussion here.



**Table 10. Comparison across six variables showing variation in JI**

		[h] ~ Ø	[s] ~ Ø	1SG	[-a] ~ [-e]	[h]~[ʔ]~Ø	N- ~ [ŋə]
Linguistic conditioning		X	X	√	X	√	X
Formal vs. informal style		X	√	√			
Social conditioning	Age	X	X	√	√	√	√
	Gender	X	X	X	X	√	√
	Ed level	X	X	X	X	√	√

We observe complex patterns of linguistic and social conditioning across the six variables, where each variable and each speaker has its own story. Age plays a role for 1SG, similar to Kurniawan’s finding, but not for C ~ Ø alternations; while neither gender nor education appear to be key factors for our three variables. The discourse effects of the 1SG variable are in line with some of Kurniawan’s findings where linguistic and social conditioning come into place. Our findings show that different variables have different linguistic conditioning because they serve different functions in the grammar, supporting Kurniawan’s suggestion that there might be such differences. This highlights the need to study variation as a system and the fact that linguistic and social systems can be interwoven.

Methodologically this work highlights the critical importance of naturalistic data and need for speaker metadata and demonstrates the fact that you cannot pick just one or two linguistic variables and assume that the social and linguistic factors that condition them will be representative of the whole system. Work looking at linguistic and social conditioning of multiple variables within a linguistic community contribute to developing a more nuanced understanding of both the linguistic and communicative competence of speakers as part of a linguistic community and a better understanding of how these work in tandem. This is all the more interesting in an emerging linguistic variety such as JI, since it is expected to exhibit more variation, at least in early stages, as speakers may have different primary languages and complex multilingual repertoires. As the variety becomes spoken more as a primary language, we might expect some of this variation to level out.

Our goal in this paper was to examine the occurrence and distribution of co-occurring variables in Jakarta Indonesian, and in doing so to underscore 1) the value of using naturalistic corpus data for these purposes; and 2) the theoretical importance of linking individual speaker practices (intraspeaker variation) with community-level practices (interspeaker variation). In the variationist sociolinguistic literature there is an ongoing tension between the almost axiomatic expectation that different social groups within a speech community will behave similarly for all variables (indeed this is the definition of a speech community, according to Labov 2006 [1966]), and the “third wave” perspective that speakers do not necessarily adhere to community patterns, but rather draw from existing linguistic repertoires in the agentive process of the linguistic construction of identity (Eckert 2012). Empirical evidence for testing the extent to which either of these hold for a single community comes from the examination of what is sometimes called co-variation or coherence (see Guy 2013; Becker 2016). In an examination of multiple variables in Brazilian Portuguese Guy argues that coherence is likely to be weaker than is often assumed, and in weighing the two approaches to multiple linguistic variables in New York City, Becker (2016) ultimately demonstrates that the two approaches are not

mutually exclusive. It is clear that further work is needed that examines co-variation generally, and particularly in the context of emerging and/or contact varieties like Jakarta Indonesian.

## References

- Abtahian, Maya R., Abigail C. Cohn, Aaron White & Yanti. 2019. Language ideologies and language shift scenarios in Indonesia. *New Ways of Analyzing Variation* 48, Panel on *What's so standard about standards?* University of Oregon, Eugene, 10–12 October 2019.
- Abtahian, Maya R., Abigail C. Cohn, Dwi Noverini Djenar & Rachel C. Vogel. 2021. Jakarta Indonesian first-person singular pronouns: form, function, and variation. *Asia-Pacific Language Variation* 7. 187–216.
- Abtahian, Maya R., Abigail C. Cohn & Yanti. 2022. Language labeling and ideology in Indonesia. *International Journal of the Sociology of Language*. 147–171. <https://doi.org/10.1515/ijsl-2021-0104>
- Ananta, Aris, Evi Nurvidya Arifin, M Sairi Hasbullah, Nur Budi Handayani & Agus Pramono. 2015. *Demography of Indonesia's Ethnicity*. Institute of Southeast Asian Studies.
- Anwar, Khaidir. 1980. *Indonesian: The development and use of a national language*. Yogyakarta: Gadjah Mada University Press.
- Becker, Kara. 2016. Linking community coherence, individual coherence, and bricolage.: The co-occurrence, of (r), raised BOUGHT, and raised BAD in New York City English. *Lingua* 172–173. 87–99.
- Chomsky, Noam. 1965. *Aspects of the theory of syntax*. Cambridge: MIT Press.
- Cohn, Abigail C. & Rachel C. Vogel. 2019. Variation in two patterns of word-initial deletion in Jakarta Indonesia: Insight from naturalistic data. In S. Calhoun, P. Escudero, M. Tabain & P. Warren (eds.), *Proceedings of the 19th International Congress of Phonetic Sciences, Melbourne, Australia 2019*. 38–42. Canberra, Australia: Australasian Speech Science and Technology Association Inc.
- Cohn, Abigail. C. & Margaret R. Renwick 2021. Embracing multidimensionality in phonological analysis. *The Linguistic Review* 38.1. 101–139. <https://doi.org/10.1515/tlr-2021-2060>
- Dardjowidjodjo, Soenjono. 1998. Strategies for a successful national language policy: The Indonesian case. *International Journal of the Sociology of Language* 130. 35–47.
- Djenar, Dwi Noverini, Michael C. Ewing, and Howard Manns. 2018. *Style and intersubjectivity in Youth Interaction*. Berlin: De Gruyter.
- Eberhard, David M., Gary F. Simons & Charles D. Fennig (eds.). 2022. *Ethnologue: Languages of the World*. Twenty-fifth edition. Dallas, Texas: SIL International. Online version: <http://www.ethnologue.com>
- Eckert, Penelope. 2012. Three waves of variation study: the emergence of meaning in the study of sociolinguistic variation. *Annual Review of Anthropology* 41. 87–100.

- Englebretson, Robert. 2007. Grammatical resources for social purposes: Some aspects of stancetaking in colloquial Indonesian conversation. In R. Englebretson (ed.), *Stancetaking in discourse: Subjectivity, evaluation, interaction*. Amsterdam: John Benjamins.
- Errington, Joseph J. 2014. Indonesian among Indonesia's languages. In E. Tagliacozzo (ed.), *Producing Indonesia: The State of the Field of Indonesian Studies*. Ithaca: SEAP Publications. 185–193.
- Ewing, Michael C. 2005. Colloquial Indonesian. In A. Adelaar & N. Himmelman (eds.), *The Austronesian Languages of Asia and Madagascar*. London: Routledge.
- Ewing, Michael C. 2019. Localising Person Reference among Indonesian Youth. In Zane Goebel, Deborah Cole & Howard Manns (eds.), *Contact Talk: The Discursive Organization of Contact and Boundaries*. 140–159. London: Routledge.
- Gil, David & Uri Tadmor. 2007. The MPI-EVA Jakarta Child Language Database. A joint project of the Department of Linguistics, Max Planck Institute for Evolutionary Anthropology and the Center for Language and Culture Studies, Atma Jaya Catholic University.
- Gil, David & Uri Tadmor. 2014. The MPI-EVA Betawi-Jakarta database. A joint project of the Department of Linguistics, Max Planck Institute for Evolutionary Anthropology and the Center for Language and Culture Studies, Atma Jaya Catholic University. <https://archive.mpi.nl/tla/islandora/search/Jakarta%20Indonesian?type=dismax>
- Guy, Greg. 2013. The cognitive coherence of sociolects: How do speakers handle multiple sociolinguistic variables? *Journal of Pragmatics* 52. 63–71.
- Ikranagara, Kay. 1975. *Melayu Betawi Grammar*. Ph.D. thesis. University of Hawaii, Manoa.
- Ikranagara, Kay. 1980. Melayu Betawi Grammar. *NUSA: Linguistics Studies in Indonesian and Languages in Indonesia* 9. 1-150.
- Kurniawan, Ferdinan. 2015. Nasal Assimilation in Jakarta Indonesian. Proceeding from AFLA 21: *Austronesian Formal Linguistics Association (AFLA)* 21. 149–165. Canberra: Australia National University.
- Kurniawan, Ferdinan. 2018. *Phonological Variation in Jakarta Indonesian: An Emerging Variety of Indonesian*. Ph.D. thesis. Cornell University.
- Labov, William. 2006. *The Social Stratification of English in New York City*, 2nd edition. New York: Cambridge University Press.
- Labov, William. 2010. *Principles of Linguistic change. Volume III: Cognitive and Cultural Factors*. Oxford: Wiley Blackwell.
- Lapoliwa, Hans 1981. *A Generative Approach to the Phonology of Bahasa Indonesia*. *Pacific Linguistics, D.* (34). Canberra: Pacific Linguistics.
- Manns, Howard, Deborah Cole & Zane Goebel. 2016. Indonesia and Indonesian. In Z. M. Goebel, D. Cole & H. Manns (eds.), *Margins, hubs, and peripheries in a decentralizing Indonesia*. (Tilburg Papers in Culture Studies, No. 162).
- Manns, Howard. 2014. Youth Radio and Colloquial Indonesian in Urban Java. *Indonesia and the Malay World* 42 (122). 43–61.

- Musgrave, Simon. 2014. Language shift and language maintenance in Indonesia. In Peter Sercombe & Ruanni Tupas (eds.), *Language, education and nation-building: Assimilation and shift in Southeast Asia*. 87–105. Basingstoke, Hampshire & UK: Palgrave Macmillan.
- Poedjosoedarmo, Soepomo. 1982. *Javanese Influence on Indonesian*. Materials in Languages of Indonesia, No. 7. Series D. Canberra: Pacific Linguistics.
- Ravindranath, Maya R. & Abigail C. Cohn. 2014. Can a language with millions of speakers be endangered? *Journal of the Southeast Asian Linguistics Society* 7. 64–75.
- Sneddon, James N. 2003a. *The Indonesian Language: Its History and Role in Modern Society*. Sydney: UNSW Press.
- Sneddon, James N. 2003b. Diglossia in Indonesian. *Bijdragen tot de Taal-, Land- en Volkenkunde* 159.4. 519–549. DOI: 10.1163/22134379-90003741
- Sneddon, James 2006. *Colloquial Jakartan Indonesian*. Canberra, Australia: Pacific Linguistics.
- Wallace, Steven. 1976. *Linguistic and Social Dimensions of Phonological Variation in Jakarta Malay* (Unpublished Ph.D. thesis). Cornell University. Corpus available at: <https://ecommons.cornell.edu/handle/1813/45568>
- Wouk, Fay. 1999. Dialect contact and koineization in Jakarta, Indonesia. *Language Sciences* 21. 61–86.