

1 **African American English and early literacy: A comparison of approaches to quantifying**
2 **nonmainstream dialect use**

3
4 Zachary Maher^{1,2*}

5 Michelle Erskine²

6 Arynn Byrd²

7 Jeffrey Harring³

8 Jan Edwards^{2,4}

9
10 ¹Program in Neuroscience and Cognitive Science

11 ²Department of Hearing and Speech Sciences

12 ³Department of Human Development and Quantitative Methodology

13 ⁴Maryland Language Science Center

14 University of Maryland, College Park; College Park, MD

15 *Corresponding Author

16 zach@umd.edu

17 0121 Taliaferro Hall, Chapel Drive, College Park, MD 20724

18 (301)-405-9789

19 **Conflict of Interest**

20 The authors claim no conflicts of interest.

21 **Funding Statement**

22 This research was supported in part by IES grant #R305A17013 to Jan Edwards and NSF grant
23 #1449815 to Colin Phillips.

24

Abstract

25 **Purpose:** Many studies have found a correlation between overall usage rates of nonmainstream
26 forms and reading scores, but less is known about which dialect differences are most predictive.
27 Here, we consider different methods of characterizing AAE use from existing assessments and
28 examine which methods best predict literacy achievement.

29 **Method:** Kindergarten and first grade students who speak AAE received two assessments of
30 dialect use and two assessments of decoding at the beginning and end of the school year. Item-
31 level analyses of the dialect-use assessments were used to compute measures of dialect usage: (1)
32 an overall feature rate measure based on the DELV-ST, (2) a subscore analysis of the DELV-ST
33 based on items that pattern together; (3) an alternative assessment where children repeat and
34 translate sentences; and (4) “repertoire” measures based on a categorical distinction of whether a
35 child used a particular feature of MAE.

36 **Results:** Models using feature rate measures provided better data-model fit than those with
37 repertoire measures, and baseline performance on a sentence repetition task was a positive
38 predictor of reading score at the end of the school year. For phonological subscores, change from
39 the beginning to end of the school year predicted reading at the end of the school year, while
40 baseline scores were most predictive for grammatical subscores.

41 **Conclusions:** The addition of a sentence imitation task is useful for understanding a child’s
42 dialect and anticipating potential areas for support in early literacy. We observed some support
43 for the idea that morphological dialect differences (i.e. irregular verb morphology) have a
44 particularly close tie to later literacy, but future work will be necessary to confirm this finding.

45

Introduction

46 Since the earliest sociolinguistic work on nonmainstream varieties of American English, there
47 has been considerable interest in potential educational implications of dialect differences. Many
48 correlational relationships have been observed between a variety of measures of dialect
49 difference and a variety of literacy outcomes, and multiple curricula have been proposed to
50 support emergent readers who speak nonmainstream dialects.

51 All of this work requires that researchers (1) have a framework for understanding what
52 dialect variation is and (2) operationalize this understanding with one or more measures of
53 participants' dialect use. Both of these steps are fraught with challenges and require
54 simplification that will fail to fully capture individuals' experience with linguistic variation. Such
55 simplification can affect the inferences we draw about the relationship between dialect
56 differences and early literacy and, in turn, affect the strategies we use to support emergent
57 readers.

58 Perhaps the most common measurement of children's dialect use is Part 1 of the
59 Diagnostic Evaluation of Language Variation-Screening Test (DELV-ST; Seymour, Roeper, de
60 Villiers, & de Villiers, 2003). This measure was designed to determine whether a child speaks a
61 nonmainstream dialect of American English (NMAE), so it focuses on the most reliably
62 produced nonmainstream features (primarily production of dental fricatives and subject-verb
63 agreement patterns). Because of this, the DELV-ST is not ideal for capturing children's
64 knowledge of Mainstream American English (MAE). This paper uses the DELV-ST along with a
65 different test, the Dialect Assessment Battery (DAB; Craig, 2014), which is explicitly designed
66 to see whether children can produce MAE-compatible features and if they can translate
67 nonmainstream dialect features into MAE.

68 We will begin with an overview of different approaches to characterizing dialect
69 differences, with the goal of adapting our measures to capture some of the insights from these
70 approaches. Next, we will review research on the relationship between dialect differences and
71 early literacy, which motivates the present study, with a focus on studies using the DELV-ST.
72 We will then present results from an ongoing study of speakers of African American English
73 (AAE) in kindergarten and first grade, comparing different approaches to measuring their
74 linguistic system and reporting the implications that these approaches have on predicting growth
75 in decoding scores.

76 **Background**

77 **Characterizing Variation**

78 One common approach to quantifying AAE involves the use of lists of AAE features, creating a
79 Dialect Density Measure (DDM) that corresponds to the rate of usage of such features. For
80 example, Washington and Craig (2002) used a list of 26 features of AAE that differ from MAE,
81 including zero copula (*They not finished eatin' yet*), multiple negation (*I don't remember nobody*
82 *havin' no motorcycle*), and variation in subject-verb agreement (*I knew you was gonna say that*).
83 They then calculated DDM as the number of tokens of any of these features divided by the
84 number of words in a given language sample. Despite the widespread usage of this approach
85 under the “DDM” terminology, we will use the more recent term Nonmainstream Form Density
86 (NMFD), following others such as McDonald and Oetting (2019), to emphasize the fact that
87 everyone speaks a dialect, and feature rate methods inherently involve a comparison to a
88 perceived standard, the “mainstream dialect.”

89 Variants of this approach have been used with tasks ranging from highly structured
90 sentence-elicitation tasks (e.g., Charity et al., 2004), to open-ended narrative tasks (e.g., Renn &

91 Terry, 2009; Craig et al., 2014). Multiple approaches to feature sets have been shown to be
92 highly correlated; for example, Renn and Terry (2009) found that a subset of just six features
93 provided comparable results to a 40-feature list in detecting style shifts in AAE-speaking sixth
94 graders, Oetting and McDonald (2002) found that type-based and token-based approaches
95 correctly categorized child speakers of AAE and Southern White English, and Oetting and Pruitt
96 (2005) found that streamlined feature lists can often be sufficient for characterizing participants'
97 dialect. Multiple approaches have also been used for the denominator in these calculations,
98 including number of opportunities to use a feature (in more structured tasks), number of words in
99 the sample (e.g., Washington & Craig, 2002; Horton-Ikard & Weismer, 2005), and number of
100 utterances in the sample (e.g., Craig & Washington, 2000). NMFD approaches have also been
101 used to study changes in dialect use over time. For example, Terry and Connor (2012) found a
102 decrease in NMFD between kindergarten and first grade, and Terry, Connor, Petscher, and Ross
103 Conlin (2012) found a decrease in NMFD over the course of first grade that leveled off in second
104 grade. These changes likely represent a combination of factors, including development within
105 AAE toward a more adultlike grammar, developing knowledge of MAE, and changes in style-
106 shifting (Beyer & Hudson Kam, 2012; Green, 2011).

107 While the feature rate approach is helpful for quantifying differences from MAE, it has
108 many limitations, which are discussed at length by Green (2011, Ch. 2). She argues that feature
109 lists are ill-equipped to characterize AAE as its own rule-governed system. This means, for
110 example, that groups of features might share underlying patterns, and NMFD would not
111 highlight this fact. Additionally, NMFD measures typically treat features that are consistent with
112 MAE as use of MAE, while only features of AAE that differ from MAE are counted as AAE use.
113 This might be appropriate for verbal *-s*, which is often considered to be absent from the AAE

114 grammar (Newkirk-Turner & Green, 2016; for criticism of this view, see Baugh, 1990;
115 Cleveland, & Oetting, 2013; Barriere et al., 2019). However, it is often the case that the “MAE”
116 feature is also available within the grammar of AAE. For example, “zero copula” commonly
117 appears on AAE feature lists, but it is variable, with overt copulas also being acceptable in AAE
118 (e.g., Wyatt, 1996; Roy, Oetting, & Moland, 2013; Newkirk-Turner, Oetting, & Stockman,
119 2014).

120 The feature rate approach also fails to account for insights from third wave
121 sociolinguistics (Eckert, 2008, 2012). Third wave sociolinguistics proceeds from the idea that
122 speakers have a range of sociolinguistic variables that they can deploy in different social
123 situations. These variables pattern together as styles, which give variables their social meanings,
124 and individuals use styles to express their ideologies about membership in different groups. This
125 approach was not necessarily developed with a focus on AAE, but it provides a framework to
126 appreciate the more nuanced nature of dialect variation and has been applied to more recent work
127 characterizing the language of African Americans (e.g., King, 2018). This is in contrast to early
128 work in sociolinguistics, which often sought to characterize idealized vernacular forms, such as
129 focusing on the idea of a pure speaker of AAE who does not “code-switch” into MAE (see
130 discussion in Wolfram, 2007; King, 2020).

131 As Snell (2013) argues, the more recent work in sociolinguistics on styles allows us to
132 use *repertoire* as an alternative framing for children’s knowledge of variation. She points out that
133 recent educational work tries to replace the “deficit narrative” (i.e. that nonmainstream features
134 indicate poor language skills) with a “difference narrative,” suggesting that nonmainstream
135 varieties are distinct, rule-governed systems. However, both of these narratives make the
136 assumption that discrete varieties of English exist, which is not borne out in the data. For

137 example, she finds that 9- and 10-year-old children in northeast England mix regional and
138 standard feature use within one discourse depending on how they are trying to socially position
139 themselves relative to their interlocutors.

140 While the UK dialect context is different from that of the United States, neither could be
141 characterized as strict diglossia, where there is clear separation between vernacular and standard
142 dialects (e.g., Auer, 2005). Also, it is not necessarily the case that a child with good
143 metalinguistic skills and knowledge of MAE will use MAE forms in the school setting; these
144 children also have compelling reasons to assert a Black identity using their speech (Ogbu, 1999),
145 and the use of different forms could be a means to assert social difference from the examiner, a
146 process known as divergence in Communication Accommodation Theory (Giles & Ogay, 2007).
147 Thus, it is possible that the mere presence of a particular MAE-compatible form (e.g., an overt
148 copula) within a child's repertoire is more important than the rate at which the child favors the
149 form over MAE-incompatible alternatives (e.g., a zero copula).

150 **Dialect Differences and Literacy**

151 Despite the criticisms of the NMFD approach, Van Hofwegen and Wolfram (2017)
152 argue that aggregate NMFD values are still useful, particularly when trying to track individuals'
153 changes in dialect usage over time and in large-scale, multidisciplinary studies in general. This
154 is a likely source of their popularity in research on the relationship between dialect differences
155 and literacy. Over the past two decades, a large body of work has developed to address the
156 extent to which *dialect mismatch*—the presence of linguistic differences between
157 nonmainstream dialects (e.g., AAE) and mainstream dialects (e.g., MAE)—impacts children's
158 literacy achievement (e.g., Connor & Craig, 2006; Labov, 1995; Terry & Scarborough, 2011).
159 The influence of dialect mismatch on literacy achievement spans various subcomponents of

160 reading, including decoding and reading comprehension, though the majority of this research
161 has focused on decoding. Studies vary in their commitments to models of reading, but we will
162 assume the Simple View of Reading, which posits that reading is a product of decoding and
163 linguistic comprehension (e.g., Gough & Tunmer, 1986; Hoover & Gough, 1990).

164 Research on language variation emerging from fields such as speech-language pathology,
165 education, linguistics, and sociolinguistics posits an inverse relationship between reading
166 achievement and the use of AAE. For example, Charity, Scarborough, and Griffin (2004)
167 examined the relationship between children’s facility with MAE via a sentence repetition task
168 that was designed to elicit features of AAE, and reading performance using a standardized
169 assessment of decoding (Woodcock Reading Mastery Test-Revised, Word Attack Subtest). They
170 calculated two NMFD scores corresponding to phonological and morphological features of AAE,
171 and they observed an inverse relationship between reading performance and the use of
172 nonmainstream dialect features for both the phonological and morphological measures.

173 Shade (2012) also used separate NMFD measures for phonological and morphological
174 features. She found that both measures were negatively correlated with decoding, but only
175 phonological NMFD was predictive of sight word reading. Research on the relationship between
176 dialect differences and literacy has generally not looked at finer-grained differences within the
177 broader categories of phonological and morphological variation, though multiple studies report
178 usage rates by feature (e.g., Washington & Craig, 2002; Craig, Thompson, Washington, &
179 Potter, 2003; Oetting & McDonald, 2002). For example, Oetting and McDonald (2002) found
180 that 100% of African American children used zero-marking on regular present tense verbs with
181 third person singular subjects, but only 70% used zero marking for irregular verbs. Such

182 separation of regular and irregular forms is also a longstanding finding in the acquisition
183 literature (e.g., Brown, 1973).

184 Other studies have used more traditional assessment methods to examine the relationship
185 between nonmainstream language variation and reading. Champion, Rosa-Lugo, Rivers, and
186 McCabe (2010) used the DELV-ST to identify speakers of nonmainstream English varieties and
187 to evaluate how performance on this screener related to a test of oral reading, the Gray Oral
188 Reading Test-Fourth Edition (GORT-4). Children who produced a greater number of
189 nonmainstream features had lower scores on the GORT-4. Others, such as Terry and Connor
190 (2012), found that a change in performance on the DELV-ST across two time points was
191 predictive of decoding skills. Children who decreased their use of nonmainstream features
192 between kindergarten and first grade had higher reading scores. This finding also highlights the
193 relevance of the time course of the relationship between NMAE use and changes in reading.

194 Collectively, this scholarship suggests that children who demonstrate higher
195 nonmainstream dialect density and less facility with varying their use of MAE in different
196 contexts (i.e., dialect-shifting) exhibit poorer literacy outcomes. This relationship remains true
197 for studies examining emergent literacy skills, such as decoding, and later literacy skills, such as
198 reading comprehension (Terry et al., 2016). Moreover, the established effects of dialect
199 mismatch on reading are above and beyond socioeconomic differences and race, factors that
200 were previously shown to obscure the relationship between dialect mismatch and reading
201 achievement (Bühler, von Oertzen, McBride, Stoll, & Maurer, 2018).

202 **Questions**

203 Here, we contrast different scoring approaches to assessments that target NMFD. We asked the
204 following research questions:

- 205 1. Do subsets of DELV-ST items pattern together in a way that corresponds to different
206 components of the AAE system?
- 207 2. Does nonmainstream dialect usage at the beginning of the school year, or change in
208 nonmainstream dialect usage during the school year better predict changes in decoding
209 abilities?
- 210 3. Does the rate of feature use or mere presence of an MAE form in an individual's
211 repertoire better predict decoding?
- 212 4. Are certain types of differences, as reflected in subscores, more useful at predicting
213 changes in decoding abilities? More specifically, given the close relationship between
214 phonology and decoding, are phonological differences more predictive of changes in
215 decoding than grammatical differences? Additionally, if forms like third singular verbal -
216 s are more indicative of a shift to MAE, will their usage be especially predictive of
217 differences in decoding?
- 218 5. Does the addition of a secondary sentence repetition task explain differences in children's
219 developing decoding abilities over and above what can be observed from DELV-based
220 measures?

221 **Methods**

222 **Participants**

223 The participants were 296 kindergarten and 260 first grade children from 12 elementary schools
224 in the Baltimore City Public Schools. All schools had a minimum of 89% African American
225 students (mean=96%) and more than 89% of students eligible for the National School Lunch
226 Program (mean=94%). All students were participating in a larger study designed to evaluate the
227 efficacy of a dialect-shifting curriculum (Edwards, 2019). Only students who did not have an

228 Individualized Education Program (IEP) were included in study; 14 students with IEPs were
229 tested but removed from analysis. A total of 69 students were excluded due to absence or transfer
230 at the second of the two testing points. Since model comparison values are only valid for models
231 that have been fit to the same data, an additional 8 participants were excluded due to a lack of
232 scorable items for a DELV subscore, withdrawal of assent on an assessment, or failure to
233 establish a ceiling score on a Basic Reading Cluster assessment due to experimenter error. Thus,
234 the present analysis used data from a total of 475 students (241 kindergarteners, 234 first
235 graders). At baseline, the mean age was 5:8 (S.D.=5 months) for kindergarteners and 6:8 for first
236 graders (S.D.=4 months), and at post, the mean age was 6:2 (S.D.=5 months) for kindergarteners
237 and 7:2 (S.D.=4 months) for first graders.

238 **Procedure**

239 All students were taken out of class for one hour of testing near the beginning of the school year
240 (October) and a second hour of testing near the end of the school year (April-May). There were
241 approximately six months between the first and second testing period. Students were tested
242 individually and received the following assessments: (1) Part 1 of the *Diagnostic Evaluation of*
243 *Language Variation-Screening Test* (DELV-ST, Seymour et al., 2003); (2) the *Dialect*
244 *Assessment Battery* (DAB, Craig, 2014); and (3) the Basic Reading Cluster (Word Attack and
245 Letter-Word Identification subtests) from the *Woodcock Johnson Achievement Test, 4th edition*
246 (WJIV; Schrank, Mather, & McGrew, 2014).

247 The DELV-ST Part 1 includes 15 items that are contrastive between AAE and MAE (see
248 examples in Table 2). Five of these items focus on phonological differences between the two
249 dialects (“DELV-Phon”), and the remaining ten items focus on differences in subject-verb
250 agreement between the dialects. Six of these items (“DELV-Irreg”) test the irregular subject-verb

251 agreement patterns of the verbs *have/has*, *don't/doesn't*, and *was/were*, and four of these items
252 (“DELV-3sg”) test use of regular verbal *-s* with third person subjects (e.g., *The girl sleeps*). The
253 DELV-ST provides a criterion score of “strong variation from MAE,” “some variation from
254 MAE,” or “no variation from MAE.” At baseline, 83% of participants were in the “strong
255 variation” category, 4% were in the “some variation” category, and 13% were in the “no
256 variation” category. At post, 80% were in the “strong variation” category, 4% were in the “some
257 variation” category, and 15% were in the “no variation” category.

258 The DAB is a non-standardized test that is designed to be used with *ToggleTalk*®, a
259 dialect shifting curriculum supplement for kindergarten and first grade students (Craig, 2014).
260 We administered a form of the DAB that was adapted to target one feature per item, and
261 sentences were recorded by the same set of four individuals who speak both AAE and MAE. The
262 DAB is composed of three subtests. Part 1, Elicited Imitation, assesses children’s ability to
263 repeat sentences produced in MAE. Part 2 assesses Dialect Recognition; students are asked to
264 state whether each sentence is produced in AAE (informal/home talk) or MAE (formal/school
265 talk). Part 3, Translation/Reformulation, asks children to translate sentences from AAE to MAE.
266 All three sections include 12 items: two sentences with plural forms, two sentences with past
267 tense, three sentences with a copula, three sentences which focus on subject-verb agreement (two
268 sentences with third-singular /s/ and one sentence with plural subject and “were”), and two
269 sentences with possessive /s/. Only Elicited Imitation (Part 1, “DAB-EI”) and Translation (Part
270 3, “DAB-TR”) were used in the present analysis.

271 The Basic Reading cluster of the WJIV assesses children’s ability to read words (Letter
272 Word Identification subtest) and nonwords (Word Attack subtest). This measure provides both
273 standard scores with a standardized mean of 100 and a standard deviation of 15 and *W*-scores,

274 which are linear raw scores. The Basic Reading standard and *W*-scores are the mean of Letter
275 Word Identification and Word Attack scores. We used the Basic Reading *W*-scores in our
276 analysis. Table 1 provides mean scores for all assessment measures used in the modeling.

277

278 Insert Table 1 about here

279

280 **Analysis**

281 **Dialect variation (DVAR) score.** For the frequency-based approach, we calculated a
282 *dialect variation* score (DVAR, a type of NMFD measure) from the DELV-ST. This score was
283 computed by dividing the total number of items that varied from MAE from the total number of
284 scorable items and multiplying by 100; a child who uses a nonmainstream form on every item
285 will receive a score of 100 (Terry et al., 2010; Terry & Connor, 2012; Terry et al., 2012).
286 Additionally, we calculated three DVAR subscores corresponding to phonological differences
287 (DVAR-Phon), irregular subject-verb agreement (DVAR-Irreg), and regular subject-verb
288 agreement (DVAR-3sg). These DVAR scores were used in this analysis.

289 The appropriateness of our selection of subscores is also supported by a confirmatory
290 factor analysis. Confirmatory factor analysis with oblimin rotation (i.e., allowing for correlated
291 factors) was performed using Mplus. Due to the discrete scale of the items, a mean and variance
292 adjusted weighted least squares estimation approach was used to extract factors (Liang & Yang,
293 2014). Each item was coded as a binary variable, where 1 corresponded to use of an AAE feature
294 that is not grammatical in MAE and 0 corresponded to an MAE-compatible utterance; all other
295 responses were treated as missing data. A two-factor model corresponding to phonological
296 (items 1-5) and morphological features (items 6-15) did not provide good data-model fit

297 ($X^2(91)=4027.76, p<0.001, RMSEA=0.20, CFI=0.63, SRMR=0.24$). However, satisfactory data-
298 model fit was obtained with a simple three-factor structure, where morphological features were
299 split into regular and irregular subject-verb agreement (i.e., factors corresponding to DELV-
300 Phon, DELV-Irreg, and DELV-3sg; $X^2(87)=449.11, p<0.001, RMSEA=0.063, CFI=0.97,$
301 $SRMR=0.07$). In order to facilitate comparison with DAB scores, we generated factor scores
302 from a five-factor structure, which also provided good data-model fit ($X^2(692)=1375.89,$
303 $p<0.001, RMSEA=0.03, CFI=0.96, SRMR=0.09$), with additional factors corresponding to AAE
304 use on the Elicited Imitation and Translation components of the Dialect Assessment Battery (see
305 Table 2). No model included cross-loadings. Regression analyses using factor scores were
306 qualitatively similar to those using DVAR subscores with DAB total scores, so only the results
307 from DVAR subscores and DAB total scores are reported here. Analyses with factor scores can
308 be found in S1.

310 Insert Table 2 about here

312 **DAB score.** Each item of the DAB received a score of 2, 1, or 0; where “2” corresponded
313 to MAE use, “1” corresponded to partial credit, and “0” corresponded to any other response,
314 including responses that involved a nonmainstream feature, since the assessment explicitly
315 prompts the use of MAE. Items received 2 points if the child produced the exact sentence of
316 MAE that was targeted, with credit awarded if proper names were changed in the child’s
317 utterance. For both Elicited Imitation and Translation, children received 1 point if their response
318 was grammatical in MAE but the sentence was modified. For Translation, children could also
319 receive 1 point if they produced the targeted MAE feature but another portion of the sentence

320 was changed, even if this change made the sentence ungrammatical in MAE. A total score out of
321 24 was calculated for each subsection.

322 **Repertoire.** To measure which phonological and morphological features that
323 differentiate AAE and MAE were in a child’s repertoire, we measured repertoire as a binary
324 variable where a score of 1 indicated that the child had used at least one form compatible with
325 MAE, and a score of 0 indicated that a child had not used an MAE form at the time point in
326 question. This was calculated for the three subcomponents of the DELV (DELV-Phon, DELV-
327 Irreg, and DELV-3sg), as well as the following features on DAB Translation: overt present-tense
328 copula (3 items); overt past tense marking (2 items); overt plural marking (2 items); and overt
329 possessive marking (2 items).¹

330 Repertoire values can be interpreted as a measurement of whether a child ever uses a
331 given form that is part of MAE, regardless of whether they sometimes (or even primarily) use a
332 different form. However, we should note that our scores are not derived from assessments that
333 target repertoire. The DELV-ST uses sentence completion to maximize the elicitation of
334 nonmainstream forms; it is designed to help clinicians identify if a child might speak a
335 nonmainstream variety of English. The DAB-Translation, on the other hand, explicitly prompts
336 children to use “school language.” Success in this task pre-supposes presence of the MAE-
337 compatible form in the child’s repertoire, but the task requires further metalinguistic skill and
338 choices of self-expression. Thus, for both tasks, it is possible for an MAE-compatible form to be
339 in a child’s repertoire but not be elicited; however, it would not make sense for a form to be
340 observed if it is not in a child’s repertoire. We are testing whether this measure has predictive
341 value, despite its limitations.

¹ Due to an oversight in stimulus preparation, children could provide a valid translation of the subject-verb agreement items without using verbal -s, so these items were excluded.

364 between baseline and post can be found in S2, and a correlation matrix of DVAR scores, DAB
365 scores, and repertoire scores at baseline is provided in Table 3.

366

367 Insert Table 3 about here

368

369 For each measure, scores at baseline were significantly correlated with scores at post at
370 the $\alpha=0.05$ level, and all of these correlations were positive. Additionally, each variable was
371 significantly correlated with DVAR-Composite at both baseline and post, with the following
372 exceptions: DAB-Rep-Plural at baseline with DVAR-Composite at baseline ($r=-0.08, p=0.07$)
373 and post ($r=-0.04, p=0.34$), and DAB-Rep-Past at baseline with DVAR-Composite at post ($r=-$
374 $0.08, p=0.08$). Because of this, DAB-Rep-Plural was not included in subsequent models.

375 **Changes in Dialect Measures Over Time**

376 For each of our measures, we confirmed the widely-observed trend of decreases in NMFD
377 throughout early school years (e.g., Terry et al., 2010). We used linear mixed-effects models to
378 measure change in each score. Each score was modeled separately, with fixed effects of time
379 point (fall or spring), grade level, and their interaction, as well as participant- and classroom-
380 level random intercepts and classroom-by-time point random slopes.

381 **DVAR and DAB.** For DVAR-Composite, there was a significant effect of time point
382 ($\widehat{\beta}^*=-0.18, S.E.=1.38, t(42.06)=-3.01, p=0.004$), indicating a decrease in NMFD for
383 kindergarteners over the course of the school year. There was a significant effect of Grade ($\widehat{\beta}^*=-$
384 $0.28, t(39.65)=-2.29, p=0.027$), indicating that first graders at baseline have lower DVAR scores
385 than kindergarteners at baseline. Finally, there was a significant interaction ($\widehat{\beta}^*=-0.17, t(41.8)=-$
386 $2.06, p=0.045$), indicating that the decrease in DVAR between fall and spring was more

387 pronounced for first graders, relative to kindergarteners. For the DVAR-Phon subscore, there
388 was a significant effect of grade ($\widehat{\beta}^*=-0.22$, $t(43.49)=-2.27$, $p=0.028$) and a significant time point
389 by grade level interaction ($\widehat{\beta}^*=-0.25$, $t(39.78)=-2.22$, $p=0.032$), suggesting that DVAR
390 phonology scores decrease, but only during first grade. For the DVAR irregular subscore, there
391 was a significant effect of time point ($\widehat{\beta}^*=-0.21$, $t(41.04)=-3.83$, $p<0.001$) and a marginal effect
392 of grade ($\widehat{\beta}^*=-0.25$, $t(37.52)=-1.93$, $p=0.061$), but no time point by grade interaction; this
393 indicates a significant increase in the use of *has*, *doesn't*, and *were* over the course of the school
394 year, with a potentially higher starting point in grade 1. There were no significant terms for the
395 DVAR-3sg subscore; this indicates that there was no increase in use of the third person singular
396 from kindergarten to first grade or from the beginning to the end of the school year.

397 For overall DAB-EI, there was a significant effect of time point ($\widehat{\beta}^*=0.33$, $t(44.55)=5.54$,
398 $p<0.001$) and grade ($\widehat{\beta}^*=0.46$, $t(42.86)=3.94$, $p<0.001$), but not an interaction, indicating that
399 usage of MAE in sentence repetition increases over the course of the school year and between
400 kindergarten and first grade. For overall DAB-TR, there was also a significant effect of time
401 point ($\widehat{\beta}^*=0.29$, $t(125.54)=4.2$, $p<0.001$) and grade ($\widehat{\beta}^*=0.45$, $t(64.74)=4.72$, $p<0.001$), as well
402 as an interaction ($\widehat{\beta}^*=0.32$, $t(124.68)=3.31$, $p=0.001$), indicating that children's ability to
403 translate sentences from AAE into MAE increases over the course of the school year and
404 between kindergarten and first grade, and this effect is more pronounced in first grade.

405

Insert Figure 1 about here

406

407

408 **Repertoire.** We ran mixed-effects logistic regression models, which are appropriate for

409 predicting binary-coded data, with fixed effects of grade and time point and their interaction, as

410 well as participant- and classroom-level random intercepts, using the glmer function of lme4. A
411 separate model was fit for each repertoire score. For overt copula usage, there was a significant
412 fixed effect of time point ($\widehat{\beta}^*=0.91, z=4.50, p<0.001$) and grade level ($\widehat{\beta}^*=0.90, z=3.93,$
413 $p<0.001$), indicating that children were more likely to have overt copula in their repertoire in the
414 spring relative to the fall and in first grade relative to kindergarten; for overt possessive usage,
415 there was a significant interaction term ($\widehat{\beta}^*=1.02, z=2.97, p=0.003$), meaning that overt
416 possessive was more likely to be in a child's repertoire at spring testing in grade 1, relative to
417 any other time point. No other terms were significant.

418 **Predicting Decoding from Dialect Measures**

419 We ran two sets of models to examine the relationship between NMAE use and reading scores.
420 In one set of models, we examined whether change in NMAE use across the school year was a
421 significant predictor of reading scores at the end of the school year. These models tested the
422 claim that being successful at the linguistic and metalinguistic demands inherent in learning to
423 dialect shift is associated with learning to read (e.g., Terry & Scarborough, 2011). In the second
424 set of models, we examined whether baseline NMAE scores were significant predictors of
425 reading at the end of the school year. These models tested the claim that learning to decode is
426 more difficult for children with higher rates of NMAE use, probably because of the greater
427 mismatch between their native dialect and the written form (e.g., Labov, 1995). Models did not
428 converge or had a singular fit using classroom-level random slopes for the relationship between
429 baseline (fall) scores and post (spring) scores, so we simplified our random effects structure to
430 include classroom-level random intercepts only. The intraclass correlation was 0.41 for the
431 unconditional model. Additionally, inclusion of grade level in models predicting reading did not
432 significantly improve model fit, so the term was dropped.

433 **Dialect variation scores (DELV-Screening Test).**

434 ***DVAR growth predicting decoding.*** We used linear mixed-effects regression models to
435 predict Basic Reading *W*-scores in the spring, with fixed effects of fall *W*-scores and the
436 difference between fall and spring DVAR scores. For the model with overall DVAR change,
437 there was a significant effect of fall Basic Reading *W*-scores ($\widehat{\beta}^*=0.85$, $t(376.32)=36.27$,
438 $p<0.001$), indicating that students with higher Basic Reading scores in the fall had higher Basic
439 Reading scores in the spring, and there was a significant effect of change in DVAR ($\widehat{\beta}^*=-0.08$,
440 $t(466.96)=-3.97$, $p<0.001$), indicating that, controlling for Basic Reading score in the fall,
441 children had higher Basic Reading scores in the spring as their NMFD decreased. Addition of
442 DAB-EI and DAB-TR score changes marginally improved model fit ($X^2(2)=5.43$, $p=0.066$),
443 driven by a marginal effect of DAB-TR ($\widehat{\beta}^*=0.04$, $t(445.55)=1.86$, $p=0.064$).

444 We fit a separate model using the three DVAR subscores as independent predictors.
445 There was a significant effect of baseline reading score ($\widehat{\beta}^*=0.85$, $t(373.4)=36.19$, $p<0.001$),
446 DVAR phonology subscore change ($\widehat{\beta}^*=-0.06$, $t(461.58)=-2.93$, $p=0.004$), and DVAR irregular
447 subscore change ($\widehat{\beta}^*=-0.05$, $t(456.4)=-2.54$, $p=0.011$), but not DVAR-3sg subscore change.
448 Again, addition of DAB-EI and DAB-TR score changes marginally improved model fit ($X^2(2)-$
449 5.43 , $p=0.066$), driven by a marginal effect of DAB-TR ($\widehat{\beta}^*=0.04$, $t(444.32)=1.77$, $p=0.078$).

450 ***DVAR baseline predicting decoding.*** Next, we used baseline DVAR scores instead of
451 change in DVAR to predict spring reading scores, controlling for baseline reading score, with a
452 classroom-level random intercept. For the model using DVAR-Composite baseline, there was a
453 significant effect of baseline Reading score ($\widehat{\beta}^*=0.84$, $t(406.47)=33.48$, $p<0.001$) and baseline
454 DVAR-Composite ($\widehat{\beta}^*=-0.05$, $t(466.25)=-2.11$, $p=0.036$). Addition of DAB-EI and DAB-TR
455 baseline significantly improved model fit ($X^2(2)=22.33$, $p<0.001$); in the full model, DVAR-

456 Composite was no longer significant, but DAB-EI baseline was ($\widehat{\beta}^*=0.10$, $t(46.16)=4.23$,
457 $p<0.001$). For the model using DVAR baseline subscores, no subscore was significant, though
458 the addition of DAB-EI and DAB-TR again significantly improved model fit ($X^2(2)=21.90$,
459 $p<0.001$), driven by a significant DAB-EI term ($\widehat{\beta}^*=0.10$, $t(457.96)=4.17$, $p<0.001$).

460 **Repertoire.** Given the exploratory nature of our repertoire scores, we used an
461 incremental model building procedure, starting with a null model with a fixed effect of baseline
462 reading score and a classroom-level random intercept. We then added the three repertoire values
463 derived from the DELV (DELV-Phon, DELV-3sg, DELV-Irreg) for model comparison, then
464 additionally included the two DAB-based repertoire measures (overt copula, overt possessive).

465 **Changes in repertoire predicting decoding.** To measure whether change in repertoire
466 predicts changes in reading, we modeled Basic Reading W -score in the spring with fixed effects
467 of baseline Basic Reading score and the change in each feature in the child's repertoire, where 1
468 indicated that the feature was added over the course of the year and 0 means it was not. In the
469 null model, Basic Reading score was significant ($\widehat{\beta}^*=0.85$, $t(380.5)=36.04$, $p<0.001$), but the
470 addition of DELV repertoire scores ($X^2(3)=1.88$, $p=0.598$) and DAB repertoire scores
471 ($X^2(2)=0.90$, $p=0.637$) did not improve model fit.

472 **Baseline repertoire predicting decoding.** Addition of baseline DELV repertoire scores to
473 the null model marginally improved fit ($X^2(3)=6.49$, $p=0.090$), and addition of DAB repertoire
474 scores significantly improved model fit ($X^2(2)=14.37$, $p<0.001$). In the full model, there was a
475 significant fixed effect of baseline reading score ($\widehat{\beta}^*=0.83$, $t(392.83)=33.97$, $p<0.001$), as well as
476 overt copula usage ($\widehat{\beta}^*=0.08$, $t(448.06)=3.56$, $p<0.001$), and a marginal effect of DELV-Irreg
477 ($\widehat{\beta}^*=0.04$, $t(458.93)=1.78$, $p=0.075$), but no other term. This indicates that controlling for
478 baseline Basic Reading score, children whose repertoire included overt copula in the present

479 tense and (possibly) MAE-compatible agreement on irregular verbs had significantly higher
480 Basic Reading scores in the spring.

481 **Model Comparison**

482 Akaike information criterion (AIC; Akaike, 1974) values for each model of Basic Reading score
483 are provided in Table 4. AIC values provide a measure of model fit that rewards parsimonious,
484 good data-model fit and penalizes overparameterized models (Anderson, 2008). In other words,
485 models are rewarded when predictors explain variance in the outcome measure, but they are
486 penalized for the number of predictors they use. In contrast to statistical tests that compare two
487 nested models, the AIC is a relative fit measure used descriptively, in which the model with the
488 lowest AIC value is considered the best-fitting model, and models with a difference in AIC value
489 of more than four relative to this best-fitting model are considered to have much weaker support.
490 We used AIC_C values, which correct for smaller sample sizes (Anderson, 2008). We refitted the
491 models using maximum likelihood estimation prior to the calculation of AIC_C values. In the
492 present analysis, the DVAR-Composite (plus DAB) approach results in the best-fitting model
493 when predicting decoding in the spring from dialect scores in the fall, with DVAR subscores
494 (plus DAB) also having some empirical support. This approach is also best overall. Of the
495 models that use changes in dialect scores to predict decoding, the DVAR measures are best.

497 Insert Table 4 about here

499 **Discussion**

500 Regardless of measurement type, we confirmed the widely-reported trend of decreased AAE use
501 and increased MAE use over the course of the school year and between kindergarten and first

502 grade. This was true for both grade levels, but it was more pronounced in first grade for some
503 measures. Given this initial validation of our measures, we return now to our research questions.

504 First, we found that a three-factor structure provides satisfactory model fit for Part 1 of
505 the DELV-ST, indicating three clusters of items: phonological items, items with regular subject-
506 verb agreement, and items with irregular subject-verb agreement. Irregular subject-verb
507 agreement spanned multiple verbs (*don't* and *haven't* with third person singular subjects, *was*
508 with plural subjects).

509 Second, models predicting decoding scores from baseline dialect measures provided
510 better fit than models predicting decoding scores from change in dialect measures. However, this
511 overall finding had a complex relationship with individual measures. Grammatical differences,
512 were only significant in baseline dialect models, with the exception of DVAR-Irreg subscores.
513 On the other hand, DVAR Phonology subscore (based on five items) was significant for models
514 that used change in dialect as predictors, but not for models that used baseline dialect measures.
515 This might indicate that knowledge of MAE grammar is a resource that children can draw upon
516 as they learn to read, and time with this resource is necessary for differences to be observed.
517 Phonology, on the other hand, is directly tied to decoding such that changes in one are predictive
518 of changes in the other. Further exploration of this potential distinction could inform future
519 research on literacy interventions. A curriculum that focuses on phonology might have an
520 immediate impact on decoding, while a grammatical one might require additional time before
521 effects are observed.

522 Models using repertoire scores had poorer data-model fit than models using NMFD, but
523 one of these models did yield a significant result for overt copula usage. We did not observe any
524 hypothesized differences between grammatical feature types. If verbal *-s* is not part of the AAE

525 grammar, we might predict that usage of verbal -s would be a particularly powerful indicator of
526 knowledge of MAE and would therefore predict reading outcomes. However, we did not observe
527 this. One possible explanation for this is that there was no increase in usage of verbal -s between
528 kindergarten and first grade or between the baseline and the end of the school year. This result is
529 consistent with other research showing that non-overt marking of third person singular shows
530 minimal change from kindergarten to fifth grade (Craig & Washington, 2004; Newkirk-Turner &
531 Green, 2016). Instead, usage of overt copulas was significant in the model predicting reading
532 from baseline repertoire, even though overt copulas are available in both AAE and MAE.

533 Finally, addition of the DAB did provide predictive value beyond the DELV-ST. The
534 Elicited Imitation subtest of the DAB at baseline predicted decoding in the spring, and this
535 proved to be a stronger predictor than any DELV measure when both were included in the same
536 model. This task is different from the DELV in that it uses MAE forms in the prompts,
537 representing a wider variety of features, and the task is repetition rather than filling in a blank. It
538 is unclear which of these differences was most important, but it is clear that even a brief, highly
539 structured task in addition to the DELV can be useful when the goal is to characterize the
540 language of a child with typical development during early literacy instruction. The elicited
541 imitation task of Charity et al. (2004), which used a picture book context, may have even
542 stronger predictive value than the simple sentence imitation task on the DAB. A comparison of
543 means and standard deviations from the Charity et al. task and our task shows more variability in
544 performance and less of a ceiling effect for Charity et al., suggesting that the picture book
545 context results in a more sensitive measure. We speculate that the picture book context promotes
546 deeper linguistic processing instead of reliance on verbal working memory.

547 One distinction that did emerge in this work is the difference between agreement in
548 irregular verbs and in other verbs. The relevant DELV-ST items loaded onto separate but
549 correlated factors, and the DVAR-Irreg measure was the only grammatical measure that was
550 significant in a growth model. This could be partially driven by the number of items (6 for
551 DELV-Irreg vs. 4 for DELV-3sg) leading to a reduction in measurement error for the DVAR-
552 Irreg measure. However, it is also plausible that children learn the irregular agreement patterns
553 without learning to use verbal *-s* on regular verbs. This aligns with early findings by Oetting and
554 McDonald (2002), who found that there were more NMAE-speaking children who used zero-
555 marking on regular third person singular verbs than who used nonmainstream subject-verb
556 agreement with *be* and *don't*. Given the frequency of these forms, the relationship between
557 knowledge of MAE agreement patterns for irregular verbs and reading could conceivably operate
558 in either causal direction. Stronger readers might learn these forms from their experiences with
559 texts, and knowledge of these frequent forms could lead to greater facility with decoding.

560 While we provided multiple ways of analyzing DELV-ST and DAB data, we are limited
561 to these two assessments in our dialect measures. Both provide highly structured elicitation
562 contexts, and previous work has made it clear that children's dialect usage differs depending on
563 the elicitation context (e.g., Craig et al., 2014; Renn and Terry, 2009). More open-ended
564 narrative tasks could provide useful comparison data and potentially elicit more nonmainstream
565 components of a child's repertoire, but such tasks are more difficult to administer and score on a
566 large scale. Such tasks also provide an opportunity to measure style shifting across contexts
567 within a given time point, rather than confounding changes in dialect and development.

568 Additionally, our reading measures are limited to measures of decoding. Specifically, the
569 Basic Reading composite score that we used is calculated from two subtests that evaluate the

570 reading of words and non-words in isolation. While these are appropriate reading measures for
571 children at this stage of schooling, it is possible that dialect differences play a distinct role in
572 passage reading (e.g., Terry et al., 2016). That is, grammatical differences between dialects may
573 be more important for passage comprehension than for decoding, since grammatical differences
574 such as agreement morphemes are more likely to appear in a passage than in isolated words.

575 We are also limited by our relatively homogeneous sample. By design, our participants
576 attended schools where students were predominantly African American and from low-SES
577 families, and all of these schools were part of the same district. Moreover, the majority of
578 participants (89%) showed at least some variation from MAE as measured by the DELV-ST.
579 Terry et al. (2012) found significant effects of race, school SES, and race-by-SES interaction in
580 predicting change in DVAR scores, so more research will be necessary to determine the degree
581 to which our results generalize to speakers of other nonmainstream dialects of English and to
582 other school settings. It is plausible that different types of experiences with variation would be
583 better captured by different measures, which is important to note when comparing studies that
584 use different populations of speakers.

585 Further research will be necessary to confirm these exploratory findings. This will
586 involve continued honing of our measurements of dialect differences. Future studies could elicit
587 larger numbers of tokens per feature and systematically vary the elicitation context to include
588 sentence repetition, sentence completion, and open-ended narrative. This would allow us to more
589 clearly determine which combination(s) of tasks and dialect differences are most predictive of
590 changes in measures of reading. This process should be repeated across multiple age ranges to
591 reflect children's evolving dialect usage. Factor analysis played only a limited role in the present
592 study, but it is a promising tool for future research. Ideally, future work would use structural

593 equation modeling not only for measures of dialect but also for studying the relationship between
594 those measures and reading scores in order to fully account for measurement error in the
595 measured variables (e.g., Johnson, Terry, Connor, and Thomas-Tate, 2017; Bühler et al., 2018).

596 As this research progresses, clinicians are faced with the challenging task of supporting
597 speakers of AAE despite having a relatively limited set of tools. One important step is to
598 characterize each child's linguistic repertoire. We have seen that the DELV-ST provides a useful
599 starting point, and it can be made more informative by grouping the items into phonology,
600 regular subject-verb agreement, and irregular subject-verb agreement. Though it might be ideal
601 to use open-form narrative tasks in a variety of settings, even a simple sentence repetition task
602 like the DAB Elicited Imitation can be helpful. As noted above, a sentence imitation task that
603 incorporates a storybook or picture description as part of the paradigm (e.g., Charity et al., 2004)
604 may result in greater semantic encoding of the sentences and thus elicit a representative range of
605 nonmainstream forms. More broadly, it is important to think of any measure of nonmainstream
606 form density as only a starting point for understanding a child's language and anticipating any
607 educational challenges from linguistic differences. The next step is providing targeted support.
608 Our results provide tentative support for the idea of focusing on areas of variable overlap
609 between MAE and AAE, such as overt copulas and irregular subject-verb agreement. This allows
610 children to draw upon their existing linguistic knowledge as they learn to read in a less familiar
611 dialect.

612
613
614
615
616
617
618
619

Acknowledgements

This research was supported in part by IES grant #R305A17013 to Jan Edwards and NSF grant #1449815 to Colin Phillips. We are very grateful to all of the children who participated in this study, their families who provided consent, and all of the amazing teachers and principals in the Baltimore City Public School System whom we have worked with. We also thank Ebony Terrell Shockley, Rebecca Silverman, Tatiana Thonesavanh and other members of the Learning to Talk Lab at the University of Maryland for their many contributions to this research program.

620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6), 716-723. <https://doi.org/10.1109/TAC.1974.1100705>
- Anderson, D. R. (2008). *Model based inference in the life sciences: a primer on evidence*. (Vol. 1). New York: Springer.
- Auer, P. (2005). Europe's sociolinguistic unity, or: A typology of European dialect/standard constellations. In *Perspectives on Variation: Sociolinguistic, Historical, Comparative*, (pp. 7–42).
- Barrière, I., Kresh, S., Aharodnik, K., Legendre, G., & Nazzi, T. (2019). The comprehension of 3rd person singular –s by NYC English-speaking preschoolers. In T. Ionin & Rispoli, M. (Eds.). *Selected Proceedings of the 7th Generative Approaches to Language Acquisition - North America Conference* (pp. 7-33). Philadelphia, PA: John Benjamins.
- Bates, D., Maechler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1-48. <https://doi.org/10.18637/jss.v067.i01>.
- Baugh, J. (1990). A survey of the suffix /-s/ analyses in Black English. In Edmondson, J. A., Feagin, C., & Mühlhäusler, P., (Eds.), *Development and Diversity: Language Variation across Time and Space: A Festschrift for Charles-James N. Bailey* (pp. 297–307). Dallas, Texas: The Summer Institute of Linguistics.
- Beyer, T. & Hudson Kam, C. L. (2012). First and second graders' interpretation of Standard American English morphology across varieties of English. *First Language*, 32(3), 365-384. <https://doi.org/10.1177/0142723711427618>
- Brown, R. (1973). *A first language: The early stages*. Harvard U. Press.

644 Bühler, J. C., von Oertzen, T., McBride, C. A., Stoll, S., & Maurer, U. (2018). Influence of
645 dialect use on early reading and spelling acquisition in German-speaking children in
646 Grade 1. *Journal of Cognitive Psychology*, 30(3), 336-360.
647 <https://doi.org/10.1080/20445911.2018.1444614>

648 Champion, T. B., Rosa-Lugo, L. I., Rivers, K. O., & McCabe, A. (2010). A Preliminary
649 Investigation of Second-and Fourth-Grade African American Students' Performance on
650 the Gray Oral Reading Test—Fourth Edition. *Topics in Language Disorders*, 30(2), 145-
651 153. <https://doi.org/10.1097/TLD.0b013e3181e04056>

652 Charity, A. H., Scarborough, H. S., & Griffin, D. M. (2004). Familiarity with School English in
653 African American children and its relation to early reading achievement. *Child
654 Development*, 75(5), 1340-1356. <https://doi.org/10.1111/j.1467-8624.2004.00744.x>

655 Cleveland, L. H. & Oetting, J. B. (2013). Children's marking of verbal s by nonmainstream
656 English dialect and clinical status. *American Journal of Speech-Language Pathology*,
657 22(4), 604-614. [https://doi.org/10.1044/1058-0360\(2013/12-0122\)](https://doi.org/10.1044/1058-0360(2013/12-0122))

658 Connor, C. M. & Craig, H. K. (2006). African American preschoolers' language, emergent
659 literacy skills, and use of African American English: a complex relation. *Journal of
660 Speech, Language, and Hearing Research*, 49(4), 771-792. [https://doi.org/10.1044/1092-
661 4388\(2006/055\)](https://doi.org/10.1044/1092-4388(2006/055))

662 Craig, H. K., & Washington, J. A. (2000). An assessment battery for identifying language
663 impairments in African American children. *Journal of Speech, Language, and Hearing
664 Research*, 43(2), 366-379. <https://doi.org/10.1044/jslhr.4302.366>

665 Craig, H. K., Kolenic, G. E. & Hensel, S. L. (2014). African American English speaking
666 students: A longitudinal examination of style shifting from kindergarten through second

667 grade. *Journal of Speech, Language, and Hearing Research*, 57(1), 143-157.
668 [https://doi.org/10.1044/1092-4388\(2013/12-0157\)](https://doi.org/10.1044/1092-4388(2013/12-0157))

669 Craig, H. K. (2014). *Toggle Talk®*. Sun Prairie, WI: Ventris Learning.

670 Craig, H. K. (2016). *African American English and the Achievement Gap*. New York: Routledge.

671 Eckert, P. (2008). Variation and the indexical field. *Journal of Sociolinguistics*, 12(4), 453-476.
672 <https://doi.org/10.1111/j.1467-9841.2008.00374.x>

673 Eckert, P. (2012). Three waves of variation study: The emergence of meaning in the study of
674 sociolinguistic variation. *Annual Review of Anthropology*, 41(1), 87-100.
675 <https://doi.org/10.1146/annurev-anthro-092611-145828>

676 Edwards, J. (2019). Dialect mismatch and learning to read: Research to practice. Plenary talk
677 presented at the 44th Annual Boston University Conference on Language Development,
678 Boston, MA, 7-10 November.

679 Fitzmaurice, G. M., Laird, N. M., & Ware, J. H. (2011). *Applied longitudinal analysis*. Hoboken,
680 NJ: John Wiley & Sons.

681 Giles, H. & Ogay, T. (2007). Communication Accommodation Theory. In Whaley, B. B. and
682 Samter, W., (Eds.), *Explaining communication: Contemporary theories and exemplars*,
683 (pp. 293-310). New York: Routledge. <https://doi.org/10.1002/9781118766804.wbiect056>

684 Gough, P. B., & Tunmer, W. E. (1986). Decoding, reading, and reading disability. *Remedial and
685 special education*, 7(1), 6-1. <https://doi.org/10.1177/074193258600700104>

686 Green, L. J. (2011). *Language and the African American Child*. Cambridge University Press.

687 Horton–Ikard, Ramonda, and Susan Ellis Weismer. "Distinguishing African American English
688 from developmental errors in the language production of toddlers." *Applied
689 Psycholinguistics* 26.4 (2005): 597-62. <https://doi.org/10.1017/S0142716405050320>

690 Hoover, W. A., & Gough, P. B. (1990). The simple view of reading. *Reading and writing*, 2(2),
691 127-16. <https://doi.org/10.1007/BF00401799>

692 Johnson, L., Terry, N. P., Connor, C. M., & Thomas-Tate, S. (2017). Effects of dialect awareness
693 instruction on nonmainstream American English speakers. *Reading and Writing*, 30(9),
694 2009-2038. <https://doi.org/10.1007/s11145-017-9764-y>

695 King, S. (2018). *Exploring social and linguistic diversity across African Americans from*
696 *Rochester, New York* (Doctoral dissertation, Stanford University).

697 King, S. (2020). From African American Vernacular English to African American Language:
698 Rethinking the Study of Race and Language in African Americans' Speech. *Annual*
699 *Review of Linguistics*, 6, 285-30. <https://doi.org/10.1146/annurev-linguistics-011619->
700 030556

701 Kuznetsova A, Brockhoff PB, Christensen RHB (2017). "lmerTest Package: Tests in Linear
702 Mixed Effects Models." *Journal of Statistical Software*, 82(13), 1–26.
703 <https://doi.org/10.18637/jss.v082.i13>

704 Labov, W. (1995). Can reading failure be reversed: A linguistic approach to the question. In V.
705 Gadsden & D. Wagner, (Eds.), *Literacy among African-American youth: Issues in*
706 *Learning, Teaching, and Schooling* (pp. 39-68). Cresskill, NJ: Hampton.

707 Liang, X., & Yang, Y. (2014). An evaluation of WLSMV and Bayesian methods for
708 confirmatory factor analysis with categorical indicators. *International Journal of*
709 *Quantitative Research in Education*, 2(1), 17-38.
710 <https://doi.org/10.1504/IJQRE.2014.060972>

711 McDonald, J. L., & Oetting, J. B. (2019). Nonword repetition across two dialects of English:
712 Effects of specific language impairment and nonmainstream form density. *Journal of*

713 *Speech, Language, and Hearing Research*, 62(5), 1381-1391.
714 https://doi.org/10.1044/2018_JSLHR-L-18-0253

715 Newkirk-Turner, B. L., & Green, L. (2016). Third person singular-s and event marking in child
716 African American English. *Linguistic Variation*, 16(1), 103-130.
717 <https://doi.org/10.1075/lv.16.1.05new>

718 Newkirk-Turner, B. L., Oetting, J. B., & Stockman, I. J. (2014). BE, DO, and modal auxiliaries
719 of 3-year-old African American English speakers. *Journal of Speech, Language, and*
720 *Hearing Research*, 57(4), 1383-1393. https://doi.org/10.1044/2014_JSLHR-L-13-0063

721 Oetting, J. B. & McDonald, J. L. (2002). Methods for characterizing participants'
722 nonmainstream dialect use in child language research. *Journal of Speech, Language, and*
723 *Hearing Research*, 45(3), 505-518. [https://doi.org/10.1044/1092-4388\(2002/040\)](https://doi.org/10.1044/1092-4388(2002/040))

724 Oetting, J. B., & Pruitt, S. (2005). Southern African-American English use across groups.
725 *Journal of Multilingual Communication Disorders*, 3(2), 136-144.
726 <https://doi.org/10.1080/14769670400027324>

727 Ogbu, J. U. (1999). Beyond language: Ebonics, proper English, and identity in a Black-American
728 speech community. *American Educational Research Journal*, 36(2), 147-184.
729 <https://doi.org/10.3102/00028312036002147>

730 Renn, J. & Terry, J. M. (2009). Operationalizing style: quantifying the use of style shift in the
731 speech of African American adolescents. *American Speech*, 84(4), 367-390.
732 <https://doi.org/10.1215/00031283-2009-030>

733 Roy, J., Oetting, J. B., & Moland, C. W. (2013). Linguistic constraints on children's overt
734 marking of BE by dialect and age. *Journal of Speech, Language, and Hearing Research*.
735 [https://doi.org/10.1044/1092-4388\(2012/12-0099\)](https://doi.org/10.1044/1092-4388(2012/12-0099))

736 Schrank, F.A., Mather, N., & McGrew, (2014). *Woodcock Johnson IV Tests of Achievement*,
737 NY: Houghton Mifflin Harcourt.

738 Seymour, H. N., Roeper, T. W., de Villiers, J., & de Villiers, P. A. (2003). *Diagnostic*
739 *Evaluation of Language Variation – Screening Test*. San Antonio, TX: Pearson.

740 Snell, J. (2013). Dialect, interaction and class positioning at school: From deficit to difference to
741 repertoire. *Language and Education*, 27(2), 110-128.
742 <https://doi.org/10.1080/09500782.2012.760584>

743 Terry, J. M., Hendrick, R., Evangelou, E., & Smith, R. L. (2010). Variable dialect switching
744 among African American children: Inferences about working memory. *Lingua*, 120(10),
745 2463-2475. <https://doi.org/10.1016/j.lingua.2010.04.013>

746 Terry, N. P. & Scarborough, H. S. (2011). The Phonological hypothesis as a valuable frame-
747 work for studying the relation of dialect variation to early reading skills. *Explaining*
748 *Individual Differences in Reading: Theory and Evidence* (pp. 97-117).

749 Terry, N. P. & Connor, C. M. (2012). Changing nonmainstream American English use and early
750 reading achievement from kindergarten to first grade. *American Journal of Speech-*
751 *Language Pathology*, 21, 78-86. [https://doi.org/10.1044/1058-0360\(2011/10-0093\)](https://doi.org/10.1044/1058-0360(2011/10-0093))

752 Terry, N. P., Connor, C. M., Thomas-Tate, S., & Love, M. (2010). Examining relationships
753 among dialect variation, literacy skills, and school context in first grade. *Journal of*
754 *Speech, Language, and Hearing Research*, 53, 126-146. [https://doi.org/10.1044/1092-](https://doi.org/10.1044/1092-4388(2009/08-0058))
755 [4388\(2009/08-0058\)](https://doi.org/10.1044/1092-4388(2009/08-0058))

756 Terry, N. P., Connor, C. M., Petscher, Y., & Ross Conlin, C. (2012). Dialect variation and
757 reading: is change in nonmainstream American English use related to reading

758 achievement in first and second grades? *Journal of Speech, Language, and Hearing*
759 *Research*, 55(1), 55-69. [https://doi.org/10.1044/1092-4388\(2011/09-0257\)](https://doi.org/10.1044/1092-4388(2011/09-0257))

760 Terry, N. P., Connor, C. M. D., Johnson, L., Stuckey, A., & Tani, N. (2016). Dialect variation,
761 dialect-shifting, and reading comprehension in second grade. *Reading and Writing*, 29(2),
762 267-295. <https://doi.org/10.1007/s11145-015-9593-9>

763 Van Hofwegen, J., & Wolfram, W. (2017). On the utility of composite indices in longitudinal
764 language study: the case of African American language. In *Panel Studies of Variation*
765 *and Change* (pp. 89-114). Routledge.

766 Washington, J. A. and Craig, H. K. (2002). Morphosyntactic forms of African American English
767 used by young children and their caregivers. *Applied Psycholinguistics*, 23(2), 209-231.
768 <https://doi.org/10.1017/S0142716402002035>

769 Wolfram, W. (2007). Sociolinguistic Folklore in the Study of African American English.
770 *Language and Linguistics Compass*, 1(4), 292-313. <https://doi.org/10.1111/j.1749->
771 [818X.2007.00016.x](https://doi.org/10.1111/j.1749-818X.2007.00016.x)

772 Wyatt, T. The acquisition of African American English copula. In: Kamhi, A.; Pollock, K.;
773 Harris, J., editors. *Communication development and disorders in African American*
774 *children*. Paul H. Brookes Publishing Co; Baltimore, MD: 1996. p. 95-116.

775

776

Supplemental Materials

777

S1: Summary table of fixed effects and lme4 model specification for each model reported in the
778 text. Models with factor score predictors are also included.

779

S2: Correlations between each pair of dialect measures at both baseline and post

780

Table 1. Means (standard deviations in parentheses) for assessment measures.

Measure Type	Measure	Kindergarten		First Grade	
		Baseline (Fall)	Post (Spring)	Baseline (Fall)	Post (Spring)
Woodcock Johnson IV	LetterWordID SS ¹	89.57 (14.42)	92.25 (14.57)	85.95 (15.32)	89.01 (17.14)
	LetterWordID <i>W</i> -Score ²	366.87 (3.23)	394.05 (29.61)	40.80 (32.20)	426.01 (34.99)
	WordAttack SS ¹	91.01 (15.38)	96.18 (15.74)	92.14 (16.56)	96.47 (16.95)
	WordAttack <i>W</i> -Score ²	42.88 (24.29)	444.08 (23.54)	45.31 (24.16)	466.15 (22.35)
	Reading SS ¹	9.39 (14.37)	94.22 (14.59)	88.96 (15.16)	92.74 (16.49)
	Reading <i>W</i> -Score ²	393.87 (26.01)	419.07 (25.40)	425.56 (26.86)	446.08 (27.73)
DVAR	Composite	82.61 (19.90)	78.63 (22.07)	75.98 (22.91)	67.84 (25.87)
	Phon	85.02 (21.42)	83.03 (24.01)	79.23 (26.05)	7.77 (32.02)
	3SG	87.52 (25.45)	86.20 (25.36)	82.73 (28.57)	78.88 (32.79)
	Irreg	76.61 (3.48)	69.55 (34.33)	67.69 (33.97)	56.12 (35.76)
DAB	Elicited Imitation (EI)	16.90 (4.61)	18.31 (4.35)	18.85 (3.86)	19.98 (3.47)
	Translation (TR)	6.49 (4.01)	7.98 (4.65)	8.75 (4.79)	11.87 (5.45)
Repertoire	DELV-Phon	.43 (.50)	.44 (.50)	.52 (.50)	.63 (.48)
	DAB-Copula	.32 (.47)	.52 (.50)	.51 (.50)	.68 (.47)
	DELV-3sg	.25 (.44)	.29 (.46)	.34 (.48)	.38 (.49)
	DELV-Irreg	.51 (.50)	.57 (.50)	.62 (.49)	.74 (.44)
	DAB-Past	.50 (.50)	.51 (.50)	.54 (.50)	.53 (.50)
	DAB-Plural	.30 (.46)	.40 (.49)	.38 (.49)	.56 (.50)
	DAB-Possessive	.19 (.39)	.17 (.37)	.18 (.38)	.32 (.47)

781

782 ¹Standardized mean is 100 and one standard deviation is 15. ²A score of 500 represents

783 normative mean achievement of a ten year old, and one standard deviation is 15.

784

785

786 Table 2. Summary of subscores based on factors used in confirmatory factor analysis

Factor	Construct	Example Item
DELV-Phon	Usage of AAE phonology	<i>smooth</i> pronounced /smuv/
DELV-Irreg	Leveling of subject-verb agreement with <i>have</i> , <i>don't</i> , and <i>was</i>	<i>The girl have a big kite.</i> <i>This girl don't like to swim.</i> <i>They was sick.</i>
DELV-3SG	Zero-marking of regular verbs with 3rd person singular subjects	<i>The boy always ride a bike.</i>
DAB-EI	Usage of AAE in a sentence repetition task	Prompt: <i>She is on the playground</i> Response: <i>She on the playground.</i>
DAB-TR	Usage of AAE when translating sentences from AAE to MAE	Prompt: <i>The boys was running.</i> Translation Goal: <i>The boys were running.</i> Response: <i>The boys was running.</i>

787

788

789 Table 3. Correlations (*r* values) among DVAR, DAB, and repertoire measures at baseline. For
 790 DVAR measures, higher values indicate greater usage of NMAE, and for DAB and Repertoire,
 791 higher values indicate greater usage of MAE. More correlation information is available in S2.

792

		DVAR				DAB		Repertoire				
		Composite	3SG	Irreg.	Phon.	EI	TR	Copula	Past	Poss.	Irreg.	Phon.
DVAR	3SG	.77***										
	Irreg.	.87***	.60***									
	Phon.	.62***	.22***	.27***								
DAB	EI	-.36***	-.28***	-.34***	-.17***							
	TR	-.39***	-.30***	-.35***	-.25***	.31***						
Rep.	Cop.	-.20***	-.13**	-.15***	-.17***	.14**	.54***					
	Past	-.10*	-.06	-.10*	-.05	.18***	.31***	-.01				
	Poss.	-.19***	-.17***	-.19***	-.09	.14**	.38***	.09	.18***			
	Irreg.	-.60***	-.34***	-.75***	-.16***	.28***	.24***	.17***	.06	.10*		
	Phon.	-.47***	-.17***	-.19***	-.79***	.16***	.15**	.15***	.00	.06	.07	
	3SG	-.65***	-.84***	-.50***	-.20***	.23***	.24***	.11*	.05	.15**	.35***	.16***

793

794 ****p*<0.001 ***p*<0.01 **p*<0.05

795

796

797 Table 4. Comparison of model fits (lower AIC_C indicates better fit; Δ_i is the difference from the
 798 best-fitting model).

	Model	df	AIC_C	Δ_i
<i>Dialect Change</i>	DVAR (Composite)	7	3826.16	0
	DVAR (Subscores)	9	3828.71	2.55
	Repertoire	9	3848.6	22.43
<i>Dialect Baseline</i>	DVAR (Composite)	7	3820.44	0
	DVAR (Subscores)	9	3824.11	3.68
	Repertoire	9	3830.52	10.08

799

800

Figure Captions

801

Figure 1: Change in dialect usage from the beginning to the end of the school year for

802

kindergarten and first grade students. Error bars represent 95% confidence intervals, and

803

small, semi-transparent points represent individual data points. (a) DVAR scores (higher=greater

804

nonmainstream form density) and DAB total scores (higher=greater MAE use; scores out of 24

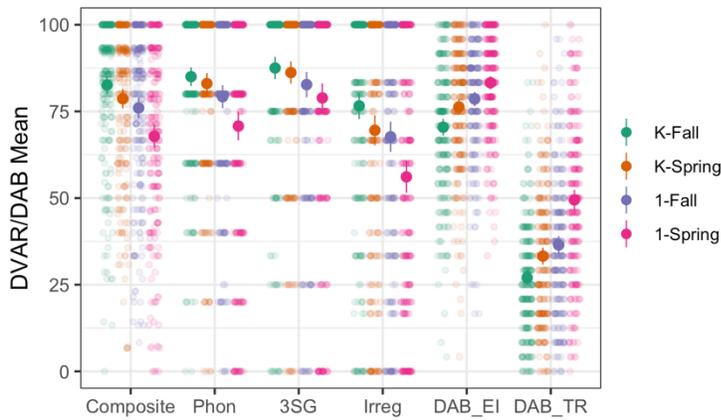
805

have been converted to percentages); (b) Repertoire Scores (of MAE-compatible feature)

806

807 (a)

808



(b)

