1

**Analyzing the Impact of Four Cognitive Constructs on NVIQ Test Performance:**

**Implications for Children with Neurodevelopmental Disorders**

Hope Sparks Lancaster [1], Erin Smolak [2], Alice Milne [3,4], Katherine R. Gordon [1],

Samantha N. Emerson [5], Claire Selin [1]

[1] Center for Childhood Deafness Language and Learning, Boys Town National Research

Hospital, Omaha NE

[2] Department of Communication Sciences and Disorders, University of South Carolina,

Columbia SC

[3] Department of Psychology, Lancaster University, Lancaster, UK

[4] Ear Institute, University College London, London, UK

[5] Training, Learning, and Readiness, Aptima, Inc., Woburn MA

Corresponding author: Hope Sparks Lancaster; 555 N 30th St., Omaha, NE, 68131

531-355-5086; hope.lancaster@boystown.org

Hope Sparks Lancaster: https://orcid.org/0000-0001-9924-9508

Erin Smolak: https://orcid.org/0000-0002-4477-7327

Alice Milne: https://orcid.org/0000-0002-2000-2842

Katherine R. Gordon: https://orcid.org/ 0000-0002-5143-8438

Samantha N. Emerson: https://orcid.org/0000-0001-7704-5934

Claire Selin: https://orcid.org/0000-0002-8392-5708

*Key words*: nonverbal intelligence, developmental language disorder, ADHD

**Abstract**

**Purpose**: Children with neurodevelopmental disorders historically exhibit lower and more variable nonverbal intelligence (NVIQ) scores compared to their typically developing peers. We hypothesize that the intrinsic characteristics of the tests themselves, particularly the cognitive constructs they assess, may account for both the lower scores and variability across tests and over time. Using a qualitative content analysis approach, we examined the extent to which key cognitive constructs are engaged in NVIQ tests and how these constructs compare across different tests.

**Methods**: Current editions of seven NVIQ tests were selected based on their relevance in clinical and research settings. Qualitative coding of constructs was developed iteratively by speech-language pathologists (SLPs) and researchers. The codes focused on cognitive domains most affected in highly prevalent neurodevelopmental conditions, including attention, receptive language, statistical learning, and working memory.

**Results**: We identified multiple sub-features for our constructs of interest. Using this coding framework, we found that NVIQ tests qualitatively differ in the extent to which these four constructs influence test performance.

**Conclusions**: Our findings suggest that understanding the impact of cognitive constructs on NVIQ tests can help explain why children with neurodevelopmental disorders exhibit lower and more unstable NVIQ scores compared to their peers. We provide recommendations for the use of NVIQ tests with neurodevelopmental disorder populations and encourage researchers and clinicians in speech and hearing sciences and psychology to use our results to inform test interpretation and selection.

Keywords: nonverbal intelligence, ADHD, DLD

51      **Analyzing the Impact of Four Cognitive Constructs on NVIQ Test Performance:**

52              **Implications for Children with Neurodevelopmental Disorders**

53              Children with neurodevelopmental disorders often score lower and show more variability

54      in performance on nonverbal intelligence tests (NVIQ) than their neurotypical peers (Gallinat &

55      Spaulding, 2014; Plante, 1998). Variability in test performance between individuals and within

56      individuals across tests and over time is particularly notable in conditions like developmental

57      language disorder (DLD) (Botting, 2005; Cole et al., 1994; Krassowski & Plante, 1997; Miller &

58      Gilbert, 2008). While past researchers have argued that lower and unstable NVIQ scores in

59      children with neurodevelopmental disorders are part of the phenotypes, we argue that the test

60      constructs artificially deflate nonverbal IQ scores due to demands on sustained attention,

61      language capacity, or familiarity with test items. For example, verbal skills could support

62      performance on NVIQ tasks; indeed, Durant et al. (2019) noted the impact of verbal skills on

63      NVIQ scores in bilingual children, albeit unequally across different tests. Furthermore, NVIQ

64      tests that use familiar objects as test items (e.g., apples), inherently tied to language, may

65      inadvertently advantage children with typical language skills (e.g., Sapir-Whorf hypothesis,

66      Gerrig & Banaji, 1994). The outcomes of NVIQ testing are crucial as they can influence the

67      eligibility category under which a child receives services and thus the types of services available

68      to them. Therefore, understanding the cognitive constructs assessed by these tests is essential.

69      Our study aims to identify the degree to which key cognitive constructs relevant to

70      neurodevelopmental disorders are required on NVIQ tests. Guided by our team's expertise in

71      attention and language impairments, these results will help improve the interpretability and

72      applicability of NVIQ test results.

73              The stated purpose of NVIQ assessments is to measure intelligence without relying on

74      language processing. Currently, NVIQ tests serve various purposes across different settings. In

75      clinical contexts, these tests help determine under which special education category to provide

76      services and assist in designing treatment plans by highlighting a child's strengths and needs. In

77    schools, Individualized Education Program (IEP) multidisciplinary teams, including educational

78    psychologists and speech-language pathologists (SLPs), use NVIQ test results for eligibility

79    verification decisions, goal setting, intervention planning, and family counseling. For example, in

80    Nebraska, intellectual ability is one of four areas considered when determining if a child qualifies

81    for services under Speech/Language Impairment (92 NAC 51.006) or a different eligibility

82    category (e.g., Intellectual Disability 006.04G; Nebraska Department of Education, Office of

83    Special Education, 2021). In research, NVIQ tests are frequently used to determine study

84    eligibility, classify children, and characterize samples (cf. Ebert & Lee, 2024; Gallinat &

85    Spaulding, 2014). Researchers also use NVIQ scores as covariates in regression analyses to

86    account for variance related to cognitive abilities (cf. Dennis et al., 2009; Elbert & Lee, 2024).

87    Thus, the selection of NVIQ tests can influence research outcomes, especially if studies differ in

88    using NVIQ as inclusion/exclusion criteria for children with specific neurodevelopmental profiles.

89    This, in turn, influences how we develop theories and understand neurodevelopmental

90    disorders.

91          Regardless of the application, it is imperative to understand exactly what NVIQ tests

92    measure. Recent discussions have highlighted the need for precision in identifying these

93    constructs (see Strand et al., 2020, p. 176). This precision is vital as NVIQ tests serve as

94    proxies for the cognitive constructs they aim to measure. Test manuals offer guidelines on

95    administration, psychometric properties, and development of NVIQ tests. These manuals detail

96    the various cognitive abilities the test developers sought to measure, often within a specified

97    theoretical framework. For example, Wechsler tests are based primarily on the Cattell-Horn-

98    Carroll theoretical framework (e.g., Wechsler, 2012, p. 1). Other tests, like the Test of Nonverbal

99    Intelligence (Brown et al., 2010), do not specify any constructs, implying that these tests are

100   "construct free." However, no NVIQ test is construct free. Typically, test manuals discuss the

101   impact of specific constructs only when a task directly aims to assess them, such as the impact

102   of working memory on tasks like digit span and picture span on Wechsler tests, but they do not

103    discuss the potential impact of other constructs such as statistical learning (Schapiro & Turk-

104    Browne, 2015).

105         Test reviews can also be used to understand the cognitive constructs on NVIQ tests

106    while comparing between tests. For example, DeThorne and Schaefer (2004) compared 16

107    NVIQ tests on administrative details, psychometric properties, and cognitive abilities tested

108    within the Cattell-Horn-Carroll framework. This valuable resource has helped many educational

109    psychologists and SLPs better understand NVIQ tests. However, to our knowledge, there are no

110    resources that compare NVIQ tests based on key constructs relevant to neurodevelopmental

111    disorder profiles (e.g., attention). Clinicians and researchers working with children with

112    neurodevelopmental disorders may be interested in understanding the potential impact of

113    cognitive constructs such as attention, receptive language, statistical learning, and working

114    memory.

115         There is a complex relationship between neurodevelopmental profiles, such as

116    developmental language disorder (DLD) and attention deficit hyperactivity disorder (ADHD), and

117    the cognitive constructs of attention, receptive language knowledge, statistical learning, and

118    working memory (e.g., Blom & Boerma, 2020; Smolak et al., 2020). Critically, children with

119    neurodevelopmental disorders, as a group, often perform lower than typically developing peers

120    on tests of attention, receptive language knowledge, statistical learning, and working memory

121    skills (Alloway & Gathercole, 2006; Duinmeijer et al., 2012; Ebert & Kohnert, 2011; Smolak et

122    al., 2020; Saffran, 2018). NVIQ tests can place high demands on attention in various ways,

123    including timed tasks. Research has found that, on average, children with language needs have

124    relatively lower attention regulation skills compared to peers with typical language (Duinmeijer et

125    al., 2012; Ebert & Kohnert, 2011; Smolak et al., 2020).

126         Language knowledge, especially receptive, can also impact NVIQ testing. Research has

127    shown that performance on NVIQ tests is verbally mediated (Durant et al., 2019a). Any use of

128    verbal instructions places demands on children's receptive language. Furthermore, NVIQ tasks

129     vary in the extent to which children can utilize language knowledge and language strategies to

130     solve tasks. For children with neurodevelopmental disorders, NVIQ performance scores may be

131     suppressed because they are less able to use language-based strategies on these tasks

132     compared to typically developing peers. Statistical learning is tapped on several NVIQ tests

133     through matrix reasoning or pattern completion tasks. Deficits in statistical learning have been

134     documented in several neurodevelopmental disorders, including DLD and ADHD (cf. Saffran,

135     2018). Thus, NVIQ test scores could, once again, be artificially lowered for children with

136     neurodevelopmental disorders due to statistical learning demands.

137          Lastly, working memory is frequently tested on NVIQ tests. Wechsler tests generally

138     include at least one working memory task, although other tests may place demands on working

139     memory by requiring a child to maintain and manipulate visual information to complete the task.

140     Some research indicates that working memory capacity is reduced in children with various

141     neurodevelopmental profiles (Alloway & Gathercole, 2006). If reduced working memory capacity

142     is a feature of neurodevelopmental disorders, this limitation may result in lower NVIQ scores

143     compared to individuals without neurodevelopmental disorders. Therefore, we posit that the

144     impact of these constructs explains why the NVIQ scores of children with neurodevelopmental

145     disorders can vary substantially across tests, especially considering the historical findings that

146     neurodevelopmental groups score lower on NVIQ tests (e.g., DLD; Gallinat & Spaulding, 2014).

147     **Purpose**

148          NVIQ tests are commonly used with children who have neurodevelopmental disorders.

149     Research has often compared these children's performance to that of neurotypical peers to

150     document and theorize differences (e.g., Botting, 2005; Karmiloff-Smith, 1998; Krassowski &

151     Plante, 1997; Thomas & Karmiloff-Smith, 2002), while test reviews have primarily focused on

152     broad theoretical frameworks, standardization, and psychometric properties (cf. DeThorne &

153     Schaefer, 2004). Currently, there is no comprehensive resource that compares NVIQ tests in

154     relation to key cognitive constructs for neurodevelopmental disorders and the degree to which

155    these constructs may impact measurement. This gap impedes clinicians' and researchers'

156    ability to make nuanced NVIQ test comparisons, often resulting in test selection driven by

157    administrative factors. Furthermore, understanding these constructs is crucial for comparing

158    outcomes across research studies, particularly when NVIQ is used as an exclusionary criterion,

159    as it can influence which children are included or excluded from a study. To improve test

160    selection in clinical and research settings, a resource comparing NVIQ tests based on the

161    constructs they assess is essential.

162         We addressed this gap by conducting a qualitative content analysis to assess the

163    relative degree of attention, receptive language, statistical learning, and working memory for

164    seven frequently used NVIQ tests. This project arose from a weekly meeting focused on DLD.

165    Additionally, select members of the team were particularly interested in ADHD and DLD, either

166    due to lived experience (author HSL) and/or programmatic lines of investigation (authors AM,

167    ES, SE, KG, HSL, CS). Thus, this qualitative content analysis focused specifically on constructs

168    significantly impacted in children with ADHD or DLD. The methodology, results, and conclusions

169    were informed and interpreted within the team's research and clinical expertise in children with

170    these conditions.

171                                      **Methods**

172    **Test Selection**

173         We selected NVIQ tests based on the Gallinat and Spaulding (2014) meta-analysis and

174    input from practicing SLPs. To reflect the tools currently available to clinicians and researchers,

175    we restricted our analysis to current editions (e.g., the Wechsler Intelligence Scales for Children

176    4th edition instead of the 3rd edition) and excluded all out-of-print or discontinued tests.

177    **Developing the Coding Scheme**

178         We used an iterative process to develop our coding scheme, with final codes and

179    operational definitions provided in Table 1. Our coding team consisted of seven researchers, six

180    holding Ph.D.s in psychology or speech and hearing sciences, and one with a master's degree

181   in speech-language pathology. We adopted a combined deductive and inductive content

182   analysis approach (cf. Bengtsson, 2016; Elo & Kyngas, 2008). This approach allowed us to start

183   with pre-defined categories (deductive) while identifying new categories and codes during

184   coding development (inductive). The coding development process involved the following steps:

185   1.  Over several weekly meetings, the team discussed concerns about how receptive

186       language and statistical learning could affect performance on NVIQ tests. These

187       discussions led to a decision to qualitatively analyze seven NVIQ tests for three broad

188       constructs: receptive language, statistical learning, and working memory, which aligned

189       with the expertise of team members.

190   2.  Two coders delineated sub-features and coding schemes for these three constructs,

191       such as verbal instructions for receptive language.

192   3.  The coders developed operational definitions for each sub-feature and initially coded the

193       tests, noting effective and ineffective elements. They also identified aspects of the tests

194       not accounted for in the initial codes.

195   4.  The full research team discussed and refined the coding scheme. During these

196       discussions, the team expanded the scope of codes to include attention and revised

197       sub-features for statistical learning.

198   5.  The coders implemented the revised scheme and made further observations.

199   6.  Additional team discussions led to further revisions of operational definitions and sub-

200       features, especially for statistical learning.

201   7.  Tests were recoded using the updated scheme.

202   8.  External feedback was sought from two researchers with Ph.D.s in psychology to ensure

203       completeness and accuracy of the coding concepts. This feedback led to the inclusion of

204       more detailed codes for attention, consideration of discontinue rules, and combining sub-

205       features for statistical learning.

206      9.  The team created and revised a flow chart for statistical learning codes, shown in Figure

207           1.

208      10. During peer review, reviewers suggested changes to our qualitative codes, including

209           adding codes for manipulatives (cf. DeThorne & Schaefer, 2004) and response type.

210           Peer review also requested further clarification in the definitions for attention, statistical

211           learning, and working memory codes.

212  **Coding and Reliability**

213      All coding was conducted simultaneously by coders, enabling real-time discussion and

214  clarification, ensuring the final codes represented consensus between the two coders, resulting

215  in 100% agreement. To synthesize the coded data, we scored each subtest, averaged these

216  scores, and assigned descriptive ranks to quantify the role of each construct within an NVIQ

217  test. The ranking descriptors used were as follows: None = < .25, Low = .25 to .49, Moderate =

218  .50 to .74, High = .75 to .99, and Very High = 1.

219                                            **Results**

220  **Description of NVIQ Tests**

221      The NVIQ tests included in our analysis were: Test of Nonverbal Intelligence - 4th edition

222  (TONI; Brown et al., 2010), Raven's Progressive Matrices (Raven & Raven, 2003), Wechsler

223  Abbreviated Scales of Intelligence (WASI-2; Wechsler, 2011), Wechsler Intelligence Scales for

224  Children (WISC-V; Wechsler, 2014), Wechsler Preschool and Primary Scales of Intelligence

225  (WPPSI-IV; Wechsler, 2012), Leiter International Performance Scale (Leiter-3; Roid et al.,

226  2013), and Kaufman Brief Intelligence Test (KBIT-2; Kaufman & Kaufman, 2004). Full citations

227  are provided in Supplemental References.

228      We compared administration details of the NVIQ tests using six general codes: number

229  of subtests, feedback provided, estimated administration time, estimated number of

230  administered items, use of manipulatives, and response format. Based on these codes, the

231  seven tests were broadly similar in administration details. Three tests had one subtest [range =

232   1 to 6] and took an estimated 20 minutes to administer [range = 15 to 40]. Individual subtests

233   ranged from 12 to 75 potential items (Supplemental Table S1). All seven tests provided

234   feedback on practice items. Four of the seven tests used manipulatives (e.g., blocks, foam

235   tiles), requiring action-based responses (e.g., placing blocks on a table). Additionally, six tests

236   allowed children to respond either nonverbally (e.g., pointing) or verbally (e.g., labeling). Table 2

237   summarizes descriptive test information and aggregate scores by coded constructs.

238   **Qualitative Coding**

239         We identified four relevant cognitive constructs from the literature on

240   neurodevelopmental disorders: attention, receptive language knowledge, statistical learning,

241   and working memory. Compared to children with typical development, children with DLD and

242   ADHD exhibit both lower overall performance and high variability in these skills (Alloway &

243   Gathercole, 2006; Smolak et al., 2020; Saffran, 2018). For each construct, we coded two or

244   three sub-features (Table 1). Across constructs, the NVIQ tests varied in the presence and

245   relative degree of these sub-features (none to very high). Aggregate scores across subtests are

246   shown in Table 2, with detailed scores for each subtest provided in Supplemental Table S1.

247         **Attention.** The system of attention is theorized to consist of three networks: alerting

248   (arousal to stimuli or vigilance), orienting (aligning to the source of sensory input), and executive

249   attention (monitoring and resolving conflict) (Petersen & Posner, 2012; Posner & Rothbart,

250   2007). Based on this theoretical framework, we coded three sub-features: penalizing changes in

251   phasic attention (i.e., "phasing" in and out; yes/no), timed responses (yes/no), and whether

252   sustained attention was required (yes/no).

253         Alertness is a state of vigilance and preparation during task performance. Tonic

254   alertness requires sustained vigilance to task goals, while phasic changes in arousal/alertness

255   moment-to-moment can negatively impact task performance (Esterman & Rothlein, 2019;

256   Petersen & Posner, 2012). For example, fluctuations in attention could cause a child to

257   prematurely reach ceiling using a consecutive discontinue rule, whereas a cumulative

258    discontinue rule would allow for phasic changes in alertness. Similarly, performance on timed

259    tasks would be more negatively impacted by phasic changes in alertness compared to non-

260    timed tasks. Finally, sustained attention requires vigilance over an extended period of time and

261    may involve aspects of both tonic alerting and orienting (Tang et al., 2015). Vigilance

262    decrements over time (due to disengagement or depletion of attentional resources) can result in

263    performance deficits.

264        The NVIQ tests ranged from little to no attention demands to moderate attention

265    demands. Only the Wechsler tests (WASI, WISC-V, WPPSI-IV) required timed responses for at

266    least one subtest. Most NVIQ tests did not require sustained attention. Tests generally allowed

267    children the opportunity to take breaks between subtests, as well as for administrators to

268    redirect the child back to the task at hand. In other words, children were not penalized for a lack

269    of sustained attention. When a test did require sustained attention, it was typically for one or two

270    subtests (e.g., WISC-V Figure Weights).

271        Most tests penalized children for phasing in and out. Four tests used multiple

272    consecutive failures as their discontinue rule. Therefore, if a child became briefly inattentive,

273    they could prematurely reach ceiling and obtain a score that underestimated their abilities for

274    the explicit constructs measured on individual subtests. The impact of phasic attention was

275    clearest for the Wechsler tests, which had moderate attention demands. These results were

276    driven by two factors: (1) the high number of subtests and (2) the stopping/discontinue rules for

277    these subtests. As demands on phasic attention accumulated over multiple subtests, children

278    with poor attention regulation were more likely to have final scores that underestimated their

279    abilities. For assessments with fewer subtests or alternative stopping rules, the impact of

280    attention performance was less pronounced.

281        **Receptive language knowledge**. We defined receptive language as the language

282    knowledge stored in a child's long-term memory that supports comprehension of verbal

283    information, whether provided by another (e.g., verbal instructions) or by the child's own internal

284     processes (e.g., understanding sub-vocal thoughts). Therefore, we identified two sub-features

285     for receptive language: verbal instructions and possible verbal strategy use.

286             All NVIQ tests required some level of receptive language (aggregate score range: 0.5–

287     0.83), and five tests had high demands. Three distinct patterns emerged across the tests:

288        1.   Moderate verbal instructions plus verbal strategy use (TONI).

289        2.   High verbal instructions with little opportunity for verbal strategy use (Raven's Matrix).

290        3.   No verbal instructions with high potential for verbal strategy use (Leiter-3).

291     For example, the TONI had an aggregate receptive language score of 0.75 (high), driven by its

292     moderate use of verbal instructions (i.e., less than 50 words and no complex syntax; see Table

293     2 and Supplemental Table S1) and a high verbal strategy score ("Which one of these goes in

294     this box?" page 5, Brown et al., 2010). Although Raven's Matrix included verbal instructions with

295     subordinate clauses, the highly abstract items reduced the likelihood of using verbal strategies,

296     as children may lack the words to describe the items. In contrast, the Leiter-3 explicitly

297     instructed administrators not to use language, instead providing suggested gestures to

298     demonstrate instructions. Additionally, the Leiter-3 often used familiar items, particularly for

299     younger children. The use of familiar items could increase the likelihood of a child employing

300     verbal strategies to solve problems, as they can rely on verbal labels stored in long-term

301     memory (e.g., "apple") to support performance.

302             We also observed that within individual NVIQ tests, receptive language demands varied

303     across subtests. This pattern is most evident in the WISC-V, which has an aggregate receptive

304     language score of 0.83 (high). This high score was primarily influenced by two subtests (visual

305     puzzles and coding), where both verbal instructions and verbal strategy use were scored as

306     high. These subtests featured complex verbal instructions and a high potential for children to

307     use verbal strategies.

308             **Statistical learning**. Statistical learning coding followed the flow diagram in Figure 1.

309     Our operational definition was based on Frost, Armstrong, and Christiansen's (2019) definition,

310    which states that statistical learning involves "perceiving and learning any forms of patterning in

311    the environment that are either spatial or temporal in nature" (p. 1130). Accordingly, we

312    identified three sub-features: pattern learning, implicitness, and cross-trial learning.

313          We first determined whether the subtest contained a pattern learning component and

314    then rated the sub-features of implicitness and cross-trial learning to evaluate how much a

315    child's ability to respond to regularities influenced their performance. By definition, patterning

316    requires more than one stimulus (an independent stimulus is not a pattern) and more than a

317    single occurrence of events in the stream (a single appearance is not a pattern) (Frost et al.,

318    2019, p. 1130).

319          The inclusion of implicitness as a sub-feature was informed by literature suggesting that

320    statistical learning and implicit learning "reflect a type of incidental pattern learning (i.e., learning

321    occurring without intention or instruction)" (Conway, 2020, p. 280) and that "[...] statistical

322    learning can occur largely automatically, without intent, without conscious awareness, and that it

323    is often implicit and incidental" (Frost et al., 2019, p. 1145). In our coding, implicitness evaluated

324    whether test instructions explicitly alerted children to the presence of a pattern.

325          The cross-trial learning sub-feature was derived from procedural learning literature,

326    which describes how "learning occurs on an ongoing basis during multiple trials" (Ullman &

327    Pierpont, 2005, p. 401).

328          All NVIQ tests included at least one subtest that tapped statistical learning. Regarding

329    pattern learning components, requirements differed more between subtests within a given NVIQ

330    test than between the tests overall. Some subtests did not require pattern learning at all (e.g.,

331    block design), while others heavily relied on it (e.g., matrix reasoning). This variability often

332    resulted in intermediate pattern learning requirements at the test level. For example:

333    •    The **WISC-V** included two subtests with high (matrix reasoning) or very high (coding)

334         pattern learning demands, while the other four subtests did not involve pattern learning.

335        As a result, the WISC-V has a low overall pattern learning demand, assuming all six

336        subtests are administered.

337    •    The **Leiter-3** included two subtests with no pattern learning demands (figure ground and

338         form completion) and two subtests with moderate demands (classification and sequential

339         order), resulting in a low overall pattern learning score.

340    These examples demonstrate the importance of understanding the influence of cognitive

341    constructs at the subtest level as well as at the overall test level.

342        There was variability in how explicit the instructions were regarding the presence of

343    patterns in subtests requiring pattern learning. The WISC-V coding instructions indicate that

344    each item will be repeated twice (explicit) but do not allow the assessor to highlight co-

345    occurrences (e.g., triangle goes with circle). The TONI and Leiter-3 provide suggested

346    instructions to indicate to the child that a pattern is expected, whereas the minimalistic

347    instructions on the WPPSI-IV do not allow explicit guidance about patterns.

348        For cross-trial learning, five of the NVIQ tests included subtests with a limited number of

349    patterns that facilitated cross-trial learning. For instance, the **KBIT-2** matrices task contains

350    approximately three patterns, each grouped together. The first 17 trials involve pairs of known

351    objects that go together (e.g., washing machine and shirt, bathtub and person).

352    In contrast:

353    •    The **Raven's Matrix** and the **WPPSI-IV** matrix reasoning subtests involve categories of

354         matrices presented randomly throughout the subtest. This design effectively eliminates

355         opportunities for cross-trial learning in these subtests.

356        **Working memory**. While various models of working memory exist (e.g., Baddeley,

357    1986; Cowan, 2001), our analysis was agnostic to the precise mechanisms underlying working

358    memory. For this study, working memory was defined as the system by which information is

359    activated in memory, held in short-term storage, and available for processing. Accordingly, we

360    coded for two sub-features: encoding novel information and holding target information in mind.

361

362        The KBIT-2 is the only NVIQ test with no working memory demands according to our

363    coding scheme (no novel information, no target holding). By contrast, the Leiter-3, WISC-V, and

364    WPPSI-IV had moderate working memory demands, while the TONI, Raven's, and WASI-2

365    relied heavily on working memory due to their use of novel items and the need to hold targets in

366    mind to complete tasks. Only the WPPSI-IV and WISC-V explicitly stated that they were

367    designed to measure working memory; thus, for these tests, the impact of working memory on

368    measurement is intentional. None of the other NVIQ tests indicated an explicit intention to

369    assess working memory.

370        Most NVIQ subtests do not explicitly assess working memory. However, a child's ability

371    to hold the target in mind can still influence performance, even when items and answers are

372    presented simultaneously. The WISC-V picture span and WPPSI-IV picture memory subtests

373    are exceptions, as they were designed specifically to assess working memory by removing the

374    target stimulus. Despite this, these subtests do not significantly increase the overall working

375    memory demands of the WISC-V (0.63) or WPPSI-IV (0.50) compared to other NVIQ tests (e.g.,

376    TONI 0.75), likely because the other subtests in these tests have minimal working memory

377    demands.

378        While simultaneous presentation may suggest that items do not need to be held in

379    memory, children who can remember the target and then select the answer without revisiting

380    the target likely have a higher working memory capacity. This ability enables consistent

381    performance across tasks with either simultaneous or non-simultaneous presentations.

382    Conversely, children lacking this ability might perform better on tasks with simultaneous

383    presentations, where activating working memory is less critical. Most tests are designed to

384    encourage the strategic approach of keeping the target in mind while selecting the answer,

385    regardless of their outward presentation.

386        Our coding for novelty showed greater variability across NVIQ tests than target holding.

387    Most tests included subtests with novel items, such as block design tasks that featured new

388    patterns of colored squares. However, the Leiter-3 and KBIT-2 primarily used familiar items. For

389    example, in the Leiter-3, early items involve recognizable objects like balloons or trees, with

390    more challenging items requiring manipulation of parts of these familiar objects, such as circle

391    segments. Similarly, the KBIT-2's matrix task included grouped patterns, with the first 17 trials

392    consisting of pairs of commonly associated objects, like a washing machine and a shirt.

393                                            **Discussion**

394        NVIQ tests aim to measure general intellectual abilities without bias from receptive or

395    expressive language skills. However, research indicates that individuals with

396    neurodevelopmental disorders often score lower and exhibit more variability on these tests

397    compared to neurotypical individuals. Traditional interpretations of these differences attribute the

398    lower scores to the phenotype of the disorders. In contrast, this study offers an alternative

399    explanation: the cognitive constructs engaged by NVIQ tests may impact performance in

400    children with neurodevelopmental disorders.

401        Our findings build on prior research by focusing on constructs associated with multiple

402    neurodevelopmental disorders and comparing how common NVIQ tests engage these

403    constructs. Our qualitative content analysis of seven major NVIQ tests highlights two key

404    findings. First, we identified how four cognitive constructs—attention, receptive language,

405    statistical learning, and working memory—are engaged across tests. These constructs, often

406    areas of difficulty for neurodevelopmental disorder populations, include various sub-features

407    (Tables 1 and 2) that could significantly influence performance. Second, we found notable

408    differences in the degree to which these constructs and their sub-features are represented

409    across the tests, as detailed in Table 2 and Supplemental Table S1.

410        The cognitive construct demands of the NVIQ tests analyzed range from minimal to

411    significant. For example, the Raven's Progressive Matrices showed no attentional demands,

412   while the KBIT-2 had no working memory requirements within our coding framework. These

413   insights provide an alternative perspective on why children with neurodevelopmental disorders

414   often achieve lower NVIQ scores compared to their peers. While prior studies have examined

415   correlations between specific child characteristics and NVIQ outcomes, our findings emphasize

416   the role of the cognitive constructs these tests assess in influencing performance. For instance,

417   a child with developmental language disorder (DLD) may score lower on NVIQ tests due to

418   receptive language demands, potentially misrepresenting their true cognitive aptitude.

419         The emphasis on particular cognitive constructs in NVIQ tests may also explain the

420   instability of NVIQ scores over time and across different tests in children with

421   neurodevelopmental disorders. For example, consider a child with DLD assessed with the

422   Leiter-3 at age five and the WISC-V at age nine. Differences in attention, receptive language,

423   and working memory demands between these tests could lead to a significant decline in NVIQ

424   scores, potentially in the range of 10 to 30 points as documented in previous research (e.g.,

425   Plante, 1998). This example underscores the importance for clinicians and researchers to

426   understand the cognitive constructs engaged by an NVIQ test when evaluating or interpreting

427   results for children with neurodevelopmental disorders.

428         Understanding the specific constructs assessed by NVIQ tests is critical for accurately

429   representing the cognitive abilities of children with neurodevelopmental disorders. This

430   knowledge enables more precise interpretations of test results and supports appropriate test

431   selection. This study contributes an important perspective to the discussion surrounding NVIQ

432   testing, urging careful consideration of the cognitive constructs engaged by these tests to

433   ensure fairer and more accurate assessments of intellectual ability in children with

434   neurodevelopmental disorders.

435   **Implications for Clinical Practice**

436         Many standardized tests (e.g., language, general cognition) assess multiple skills,

437   including some not explicitly identified in the test manual. SLPs should be mindful of the key

438   constructs that could influence test performance. A deeper understanding of NVIQ assessments

439   enables SLPs to make more person-centered decisions regarding eligibility, diagnoses,

440   treatment plans, and family counseling. Given the high-stakes nature of NVIQ assessments,

441   SLPs can collaborate with school psychologists and neuropsychologists to select suitable,

442   appropriate tests (ASHA Assessment and Teaming, n.d.). For example, when language poses a

443   challenge, tests like the Raven's Progressive Matrices and Leiter-3 are better choices for

444   assessing nonverbal skills with minimal language interference.

445        Returning to the Nebraska state guidelines, consider a child with suspected DLD being

446   evaluated in Nebraska. As part of the assessment, the multidisciplinary IEP team administers

447   the WISC-V to measure intellectual ability. Children with DLD typically score lower on verbal

448   tasks due to their language-learning profiles. However, our qualitative findings indicate that

449   nonverbal subtests on the WISC-V have high receptive language demands. These demands

450   could lead clinicians to underestimate the child's nonverbal abilities, potentially resulting in

451   misclassifications such as Intellectual Disability. Using our findings, the SLP on the IEP team

452   could advocate for selecting a different NVIQ test for this child, leading to a more appropriate

453   assessment. Effective interprofessional collaboration in selecting evaluation tools is essential for

454   accurate data interpretation and proper eligibility classifications.

455   **Implication for Research**

456        Our findings have significant implications for research involving NVIQ tests. These tests

457   are commonly used to determine study eligibility, classify children, match participants by NVIQ,

458   characterize research samples, and control for NVIQ in regression analyses. Concerns

459   regarding the use of NVIQ tests in these contexts have been raised in previous studies (Dennis

460   et al., 2009; Earle et al., 2017; Norbury et al., 2016). Our results suggest that researchers must

461   carefully select NVIQ tests to align with their study objectives, as failure to do so can introduce

462   confounds that affect research results.

463        The impact of NVIQ test selection on research findings can be illustrated with an

464    example. Consider a study on working memory in children with DLD that uses the TONI to

465    screen participants with a cutoff score of 85, as justified by prior research (e.g., Leonard, 2014).

466    Due to the TONI's high demands on receptive language and working memory, the participants

467    included in the study would likely have average or above-average working memory skills. This

468    could lead researchers to erroneously conclude that working memory does not differ

469    significantly between children with DLD and those without, thereby affecting the study's

470    generalizability and theoretical models of working memory's role in language learning.

471    Alternatively, using tests such as the Raven's Progressive Matrices, Leiter-3, or KBIT-2 might

472    include a broader range of abilities, reducing confounds and better targeting the construct of

473    interest. This highlights the importance of aligning NVIQ test selection with study goals rather

474    than solely relying on precedent, as test choice can significantly influence study outcomes and

475    interpretations.

476    **Strengths and Limitations**

477        This study has three main limitations. First, our qualitative coding scheme may not

478    encompass every construct assessed by the NVIQ tests analyzed. We focused on identifying

479    common constructs and their sub-features, particularly those known from previous research to

480    affect NVIQ test performance in children with neurodevelopmental disorders such as ADHD and

481    DLD. However, some NVIQ tests may assess additional constructs that are not common across

482    multiple tests. For example, there is considerable overlap between NVIQ and executive

483    functioning, which we did not include in our coding scheme. We encourage clinicians and

484    researchers to examine the impact of executive functioning on NVIQ test performance in

485    children with neurodevelopment disorders.

486        Second, we did not analyze every available NVIQ test. Instead, we selected tests that

487    are currently in print and widely used in clinics, schools, and research settings. This selection

488    was informed by clinical SLP input and the tests' prevalence in research, and the chosen tests

489     have been staples in various settings for over 30 years. Our coding scheme is provided in Table

490     1, and we invite SLPs and researchers to apply it to other NVIQ tests. To facilitate this, we have

491     established an OSF project for community use of this coding framework, aiming to support the

492     selection and interpretation of NVIQ tests in clinical and research contexts.

493            Third, we chose to interpret our results with a focus on children with ADHD and DLD.

494     However, there are several other neurodevelopmental disorders that exhibit similar weaknesses

495     in attention, language, statistical learning, and working memory. Clinicians and researchers may

496     wish to reinterpret our findings to consider how these cognitive constructs could impact NVIQ

497     test performance for children with autism spectrum disorder, intellectual disability, or specific

498     learning disabilities (e.g., dyslexia).

499     **Conclusion**

500            This study explored an alternative explanation for the observed lower and more variable

501     NVIQ scores among children with neurodevelopmental disorders: the influence of attention,

502     receptive language knowledge, statistical learning, and working memory on NVIQ test

503     performance. Building on prior work, we focused on these constructs because they are

504     associated with multiple neurodevelopmental disorders. Although we did not address disorders

505     such as autism in this study, we strongly recommend that researchers with expertise in this area

506     explore the potential implications of these constructs further.

507            Our multidisciplinary team used qualitative coding to evaluate the extent to which these

508     four key constructs could affect performance on NVIQ measures in children with

509     neurodevelopmental disorders. Of these constructs, only attention and working memory were

510     discussed in test manuals, while the impact of receptive language knowledge and statistical

511     learning was not. These constructs appeared to varying degrees across all the NVIQ tests we

512     analyzed.

513            It is essential for clinicians and researchers to consider the influence of these constructs

514     when selecting and interpreting NVIQ tests. Our qualitative coding framework, used in

515     combination with test manuals and reviews (e.g., DeThorne & Schaefer, 2004), can serve as a

516     valuable resource to aid in test selection and interpretation.

517     **Acknowledgments**

518             We extend our gratitude to everyone who provided constructive feedback on this project

519     and manuscript.

520     **Data Sharing**

521     https://osf.io/3vb57/?view_only=f53f2d3f8635491fa696b38fb3664e4a

**References**

522

523    Alloway, T. P., & Gathercole, S. E. (Eds.). (2006). *Working Memory and Neurodevelopmental*

524          *Disorders*. Psychology Press. https://doi.org/10.4324/9780203013403

525    Baddeley, A. (1986). *Working memory* (pp. xi, 289). Clarendon Press/Oxford University Press.

526    Bengtsson, M. (2016). How to plan and perform a qualitative study using content

527          analysis. *NursingPlus Open*, *2*, 8-14.

528    Blom, E., & Boerma, T. (2020). Do children with developmental language disorder (DLD) have

529          difficulties with interference control, visuospatial working memory, and selective

530          attention? Developmental patterns and the role of severity and persistence of DLD.

531          *Journal of Speech, Language, and Hearing Research*, *63*(9), 3036–3050.

532          https://doi.org/10.1044/2020_JSLHR-20-00012

533    Brown, L., Sherbenou, R. J., & Johnsen, S. K. (2010). *Test of Nonverbal Intelligence, Fourth*

534          *Edition*. Pro-Ed.

535    Cole, K. N., Mills, P. E., & Kelley, D. (1994). Agreement of assessment profiles used in cognitive

536          referencing. *Language, Speech, and Hearing Services in Schools*, *25*(1), 25–31.

537          https://doi.org/10.1044/0161-1461.2501.25

538    Conway, C. M. (2020). How does the brain learn environmental structure? Ten core principles

539          for understanding the neurocognitive mechanisms of statistical learning. *Neuroscience &*

540          *Biobehavioral Reviews*, *112*, 279–299. https://doi.org/10.1016/j.neubiorev.2020.01.032

541    Cowan, N. (2001). The magical number 4 in short-term memory: A reconsideration of mental

542          storage capacity. *Behavioral and Brain Sciences*, *24*(1), 87–114.

543          https://doi.org/10.1017/S0140525X01003922

544    Dennis, M., Francis, D. J., Cirino, P. T., Schachar, R., Barnes, M. A., & Fletcher, J. M. (2009).

545          Why IQ is not a covariate in cognitive studies of neurodevelopmental disorders. *Journal*

546          *of the International Neuropsychological Society*, *15*(3), 331–343.

547          https://doi.org/10.1017/S1355617709090481

548    DeThorne, L. S., & Schaefer, B. A. (2004). A guide to child nonverbal IQ measures. *American*

549         *Journal of Speech-Language Pathology*, *13*(4), 275–290. https://doi.org/10.1044/1058-

550         0360(2004/029)

551    DeThorne, L. S., & Watkins, R. V. (2001). Listeners' perceptions of language use in children.

552         *Language, Speech, and Hearing Services in Schools*, *32*(3), 142–148.

553         https://doi.org/10.1044/0161-1461(2001/012)

554    Duinmeijer, I., de Jong, J., & Scheper, A. (2012). Narrative abilities, memory and attention in

555         children with a specific language impairment. *International Journal of Language &*

556         *Communication Disorders*, *47*(5), 542–555. https://doi.org/10.1111/j.1460-

557         6984.2012.00164.x

558    Durant, K., Peña, E., Peña, A., Bedore, L. M., & Muñoz, M. R. (2019a). Not all nonverbal tasks

559         are equally nonverbal: Comparing two tasks in bilingual kindergartners with and without

560         Developmental Language Disorder. *Journal of Speech, Language, and Hearing*

561         *Research*, *62*(9), 3462–3469. https://doi.org/10.1044/2019_JSLHR-L-18-0331

562    Durant, K., Peña, E., Peña, A., Bedore, L. M., & Muñoz, M. R. (2019b). Not all nonverbal tasks

563         are equally nonverbal: Comparing two tasks in bilingual kindergartners with and without

564         Developmental Language Disorder. *Journal of Speech, Language, and Hearing*

565         *Research*, *62*(9), 3462–3469. https://doi.org/10.1044/2019_JSLHR-L-18-0331

566    Earle, F. S., Gallinat, E. L., Grela, B. G., Lehto, A., & Spaulding, T. J. (2017). Empirical

567         implications of matching children with specific language impairment to children with

568         typical development on nonverbal IQ. *Journal of Learning Disabilities*, *50*(3), 252–260.

569         https://doi.org/10.1177/0022219415617165

570    Ebert, K. D., & Kohnert, K. (2011). Sustained attention in children with primary language

571         impairment: A meta-analysis. *Journal of Speech, Language, and Hearing Research:*

572         *JSLHR*, *54*(5), 1372–1384. https://doi.org/10.1044/1092-4388(2011/10-0231)

573    Ebert, K. D., & Lee, H. (2024). Individual predictors of language treatment response in children

574             with developmental language disorder: A systematic review. *Journal of Speech,*

575             *Language, and Hearing Research*, *67*(8), 2708–2728.

576             https://doi.org/10.1044/2024_JSLHR-23-00665

577    Ebert, K. D., Rak, D., Slawny, C. M., & Fogg, L. (2019). Attention in bilingual children with

578             developmental language disorder. *Journal of Speech, Language, and Hearing Research*,

579             *62*(4), 979–992. https://doi.org/10.1044/2018_JSLHR-L-18-0221

580    Elo, S., & Kyngäs, H. (2008). The qualitative content analysis process. *Journal of Advanced*

581             *Nursing, 62(*1), 107-115.

582    Esterman, M., & Rothlein, D. (2019). Models of sustained attention. *Current Opinion in*

583             *Psychology*, *29*, 174–180. https://doi.org/10.1016/j.copsyc.2019.03.005

584    Frost, R., Armstrong, B. C., & Christiansen, M. H. (2019). Statistical learning research: A critical

585             review and possible new directions. *Psychological Bulletin*, *145*(12), 1128–1153.

586             https://doi.org/10.1037/bul0000210

587    Gallinat, E., & Spaulding, T. J. (2014). Differences in the performance of children with specific

588             language impairment and their typically developing peers on nonverbal cognitive tests: A

589             meta-analysis. *Journal of Speech, Language, and Hearing Research: JSLHR*, *57*(4),

590             1363–1382. https://doi.org/10.1044/2014_JSLHR-L-12-0363

591    Gerrig, R. J., & Banaji, M. R. (1994). CHAPTER 8—Language and Thought. In R. J. Sternberg

592             (Ed.), *Thinking and Problem Solving* (Vol. 2, pp. 233–261). Academic Press.

593             https://doi.org/10.1016/B978-0-08-057299-4.50014-1

594    Kaufman, A. S., & Kaufman, N. L. (2004). *Kaufman Brief Intelligence Test—Second Edition*.

595             Pearson.

596    Krassowski, E., & Plante, E. (1997). IQ variability in children with SLI: Implications for use of

597             cognitive referencing in determining SLI. *Journal of Communication Disorders*, *30*(1), 1–

598             9. https://doi.org/10.1016/0021-9924(95)00052-6

599 Miller, C. A., & Gilbert, E. (2008). Comparison of performance on two nonverbal intelligence

600     tests by adolescents with and without language impairment. *Journal of Communication*

601     *Disorders*, *41*(4), 358–371. https://doi.org/10.1016/j.jcomdis.2008.02.003

602 Nebraska Department of Education, Office of Special Education. (2021). *Determining Special*

603     *Education Eligibility—Speech Language Impairment*. https://cdn.education.ne.gov/wp-

604     content/uploads/2021/01/Eligibility-Guidelines-SLI-COMBINED1.pdf

605 Norbury, C. F., Gooch, D., Wray, C., Baird, G., Charman, T., Simonoff, E., Vamvakas, G., &

606     Pickles, A. (2016). The impact of nonverbal ability on prevalence and clinical

607     presentation of language disorder: Evidence from a population study. *Journal of Child*

608     *Psychology and Psychiatry, and Allied Disciplines*, *57*(11), 1247–1257.

609     https://doi.org/10.1111/jcpp.12573

610 Petersen, S. E., & Posner, M. I. (2012). The attention system of the human brain: 20 years after.

611     *Annual Review of Neuroscience*, *35*(Volume 35, 2012), 73–89.

612     https://doi.org/10.1146/annurev-neuro-062111-150525

613 Plante, E. (1998). Criteria for SLI. *Journal of Speech, Language, and Hearing Research*, *41*(4),

614     951–957. https://doi.org/10.1044/jslhr.4104.951

615 Posner, M. I., & Rothbart, M. K. (2007). Research on attention networks as a model for the

616     integration of psychological science. *Annual Review of Psychology*, *58*(Volume 58,

617     2007), 1–23. https://doi.org/10.1146/annurev.psych.58.110405.085516

618 Raven, J., & Raven, J. (2003). Raven Progressive Matrices. In *Handbook of nonverbal*

619     *assessment* (pp. 223–237). Kluwer Academic/Plenum Publishers.

620     https://doi.org/10.1007/978-1-4615-0153-4_11

621 Roid, G. H., Mille, L. J., Pomplun, M., & Koch, C. (2013). *Leiter International Performance*

622     *Scale—Third Edition*. Stoelting Company.

623 Saffran, J. R. (2018). Statistical learning as a window into developmental disabilities. *Journal of*

624     *Neurodevelopmental Disorders*, *10*(1), 35. https://doi.org/10.1186/s11689-018-9252-y

625    Schapiro, A., & Turk-Browne, N. (2015). Statistical Learning. In *Brain Mapping* (pp. 501–506).

626         Elsevier. https://doi.org/10.1016/B978-0-12-397025-1.00276-1

627    Smolak, E., McGregor, K. K., Arbisi-Kelm, T., & Eden, N. (2020). Sustained attention in

628         developmental languag disorder and its relation to working memory and language.

629         *Journal of Speech, Language, and Hearing Research: JSLHR*, *63*(12), 4096–4108.

630         https://doi.org/10.1044/2020_JSLHR-20-00265

631    Stark, R. E., & Tallal, P. (1981). Selection of children with specific language deficits. *The Journal*

632         *of Speech and Hearing Disorders*, *46*(2), 114–122. https://doi.org/10.1044/jshd.4602.114

633    Strand, J. F., Ray, L., Dillman-Hasso, N. H., Villanueva, J., & Brown, V. A. (2020).

634         Understanding speech amid the jingle and jangle: Recommendations for improving

635         measurement practices in listening effort research. *Auditory Perception & Cognition*,

636         *3*(4), 169–188. https://doi.org/10.1080/25742442.2021.1903293

637    Swisher, L., & Plante, E. (1993). Nonverbal IQ tests reflect different relations among skills for

638         specifically language-impaired and normal children: Brief report. *Journal of*

639         *Communication Disorders*, *26*(1), 65–71. https://doi.org/10.1016/0021-9924(93)90016-4

640    Tang, Y.-Y., Hölzel, B. K., & Posner, M. I. (2015). The neuroscience of mindfulness meditation.

641         *Nature Reviews Neuroscience*, *16*(4), 213–225. https://doi.org/10.1038/nrn3916

642    Ullman, M. T., & Pierpont, E. I. (2005). Specific Language Impairment is not Specific to

643         Language: The Procedural Deficit Hypothesis. *Cortex*, *41*(3), 399–433.

644         https://doi.org/10.1016/S0010-9452(08)70276-4

645    Wechsler, D. (2011). *Wechsler Abbreviated Scale of Intelligence—Second Edition*. Pearson.

646    Wechsler, D. (2012). *Wechsler Preschool and Primary Scale of Intelligence—Fourth Edition*.

647         Pearson.

648    Wechsler, D. (2014). *Wechsler Intelligence Scale for Children—Fifth Edition*. Pearson.

649

List of Figures:

Figure 1. Flow chart for scoring the sub-features related to statistical learning.


Supplemental file descriptions:

Supplemental Table 1 presents the qualitative coding for all nonverbal intelligence tests by sub-test, along with

the aggregate score calculations.