1

2

3

4   Playing with BEARS: Balancing Effort, Accuracy, and Response Speed in a Semantic Feature

5   Verification Anomia Treatment Game.

6

7   William S. Evans,[1]* Robert Cavanaugh,[1,2] Yina Quique,[3] Emily Boss,[4] Jeffrey J. Starns,[5] and

8   William D. Hula[2,1]

9   [1]Department of Communication Sciences and Disorders, University of Pittsburgh, Pittsburgh,

10   PA, USA.

11   [2]Geriatric Research Education and Clinical Center, VA Healthcare System, Pittsburgh, PA, USA.

12   [3]Center for Education in Health Sciences, Northwestern University, Chicago, IL, USA.

13   [4]Integrative Reconnective Aphasia Therapy, Pittsburgh, PA, USA.

14   [5]Department of Psychology, University of Massachusetts, Amherst, MA, USA.

15

16   *Corresponding author: will.evans@pitt.edu

17

**Abstract:**

18

19      **Purpose:** The purpose of this study was to develop and pilot a novel treatment framework called

20      *BEARS* (Balancing Effort, Accuracy, and Response Speed). People with aphasia (PWA) have

21      been shown to maladaptively balance speed and accuracy during language tasks. BEARS is

22      designed to train PWA to balance speed-accuracy tradeoffs and improve system calibration (i.e.,

23      to adaptively match system use with its current capability), which was hypothesized to improve

24      treatment outcomes by maximizing retrieval practice and minimizing error learning. In this

25      study, BEARS was applied in the context of a semantically-oriented anomia treatment based on

26      semantic feature verification (SFV).

27      **Methods:** Nine PWA received 25 hours of treatment in a multiple baseline single-case series

28      design. BEARS + SFV combined computer-based SFV with clinician-provided BEARS meta-

29      cognitive training. Naming probe accuracy, efficiency, and proportion of "pass" responses on

30      inaccurate trials were analyzed using Bayesian generalized linear mixed-effect models.

31      Generalization to discourse and correlations between practice efficiency and treatment outcomes

32      were also assessed.

33      **Results:** Participants improved on naming probe accuracy and efficiency of treated and untreated

34      items, although untreated item gains could not be distinguished from the effects of repeated

35      exposure. There were no improvements on discourse performance, but participants demonstrated

36      improved system calibration based on their performance on inaccurate treatment trials, with an

37      increasing proportion of "pass" responses compared to paraphasia or timeout nonresponses. In

38      addition, levels of practice efficiency during treatment were positively correlated with treatment

39      outcomes, suggesting that improved practice efficiency promoted greater treatment

40      generalization and improved naming efficiency.

41     **Conclusions:** BEARS is a promising, theoretically-motivated treatment framework for

42     addressing the interplay between effort, accuracy, and processing speed in aphasia. This study

43     establishes the feasibility of BEARS + SFV and provides preliminary evidence for its efficacy.

44     This study highlights the importance of considering processing efficiency in anomia treatment, in

45     addition to performance accuracy.

**Introduction**

Aphasia is a language disorder caused by stroke and other acquired brain injuries that affects roughly one-third of stroke survivors and more than 2 million people in the United States (Simmons-Mackie, 2018). Anomia, the inability to successfully retrieve and produce words, is a cardinal feature of aphasia (Goodglass, 1980) and experienced to some degree by all people with aphasia (PWA). Therefore, it is important to continue to improve anomia treatment outcomes, and the current work attempts to contribute to this endeavor. The current study piloted a novel game-based intervention which combined an established semantically-oriented anomia treatment (Semantic Feature Verification; Kiran & Roberts, 2010) with feedback and clinician-provided "system calibration training" (described below), designed to help PWA to balance speed, accuracy, and effort during word retrieval.

In the sections to follow, we will explain our conceptualization of system calibration and adaptation deficits in aphasia and explain how they apply to speed-accuracy tradeoffs and retrieval effort in anomia rehabilitation and functional communication. This will in turn motivate our novel system calibration training framework, BEARS ("Balancing Effort, Accuracy, and Response Speed"). The introduction will conclude with goals and study predictions for the current pilot.

*System calibration and adaptation deficits in aphasia*

In their classic work on Adaptation Theory, Heeschen and colleagues (Heeschen & Schegloff, 1999; Kolk & Heeschen, 1990; Kolk & Heeschen, 1992) argued that a distinction should be made between aphasia symptoms caused by underlying impairments ('impairment symptoms') and those caused by an individual's response to those impairments ('adaptation

69    symptoms'). Their key evidence came from observations regarding the nature of 'telegraphic

70    speech' in Broca's aphasia, where individuals who typically produced single key content words

71    at a slow rate in spontaneous speech instead produced lengthier paragrammatic output when

72    directed to do so in more constrained contexts such as a sentence elicitation task (Kolk &

73    Heeschen, 1992). In contrast, individuals with Wernicke's aphasia display task insensitivity,

74    consistently produced paragrammatic output regardless of task. As a result, they argued that two

75    key symptoms of the Broca's aphasia, slowed speaking rate and telegraphic speech, were

76    adaptation symptoms, reflecting a strategic adaptation to an underlying grammatical output

77    impairment. Their key claim, that a distinction should be made between an individual's core

78    linguistic impairments and their strategic response to these impairments, applies beyond the

79    context of understanding classic aphasia syndromes and has wide-ranging implications.

80        Building upon Adaptation Theory, we have argued that language performance and

81    communication success in aphasia are determined by a combination of impairment and

82    adaptation factors, with ultimate performance based on how well an individual makes use of

83    their language system in its current state (Evans et al., 2019). PWA who do not respond well to

84    their language system changes may demonstrate *adaptation deficits*, poorer-than-necessary

85    language performance exacerbated by the use of maladaptive strategies and habitual responses

86    (e.g., an over-reliance on ineffectual self-cuing approaches, consistently struggling to retrieve

87    difficult words with considerable effort and frustration in contexts where this rarely results in

88    success).

89        PWA who make effective use of their current core language capability demonstrate good

90    system calibration, calibrating the demands they make of their system to its current capabilities

91    in ways that are most likely to result in success. One analogy we have used with PWA to

92    describe this concept is driving a car with a manual transmission. A good driver makes best use

93    of the transmission and engine in its current condition (e.g., knows how to work with a worn-out

94    clutch). However, someone unfamiliar with the car or a new driver may not know how to apply

95    the necessary finesse, and thus may experience unnecessary issues such as stalling the engine or

96    grinding the clutch. Adaptation deficits in aphasia consist of behaviors such as making repeated

97    inaccurate retrieval attempts with increasing frustration instead of moving on or switching to an

98    alternative communication strategy. A schematic for understanding the relationship between

99    system capability, use, and calibration can be seen in Figure 1.

**Core language system capability:**
- Underlying language impairment.
- Source of impairment symptoms in Adaptation Theory.
- Aspects of task performance *not* attributable to differences in strategy/ approach.

**System use and deployment:**
- How the core language system is engaged.
- Source of adaptation symptoms in Adaptation Theory.
- Includes both deliberate strategies and habitual use.
- Likely at least partially reliant on domain-general cognitive abilities.

**Adaptive system use: good system calibration**
- Allows best-possible performance.
- System used with finesse.
- Well-calibrated practice may improve treatment efficiency for both restorative and compensatory approaches.

**Maladaptive system use: poor system calibration:**
- Leads to *adaptation deficits* (unnecessarily poor performance).
- Poorly-calibrated practice may a) Reduce restorative effects of drill-based treatment, and b) Reinforce maladaptive communication habits.

**Figure 1.** Schematic representing how core system capability (impairment) and its deployment (adaptation) together determine overall language task and communication performance, leading to *adaptive* or *maladaptive system calibration*. The BEARS training framework is intended to improve system calibration specifically as they relate to the level of effort and speed-accuracy tradeoffs during language performance.

100

101          The concept of flexible adaptation to current capability has often been associated with

102    compensatory treatments and augmentative alternative communication (Hunt et al., 2002). In the

103    current work, we wish to expand on this idea to propose that adaptive system calibration also

104    makes best use of the original intended modality (e.g., successfully producing a difficult word

105    after taking a breath to relax instead of needing to shift to an alternative communication

106    strategy). We also propose that in drill-based treatment tasks, adaptive system calibration can be

107    defined as engaging the language system in ways that improve treatment outcomes. Broadly

108    construed, adaptive system calibration is about PWA making best-possible use of their current

109    system to maximize language performance during treatments or functional communication

110    activities.

111         A benefit of this conceptualization of language performance in aphasia is that treatment

112    can target PWA's underlying language impairments (e.g., strengthen specific retrieval

113    mechanisms), strategic and habitual response to these impairments, or both. For instance, recent

114    work has demonstrated that a brief 5-session mindfulness meditation intervention for PWA

115    temporarily improves verbal fluency (Marshall et al., 2018), even such training is unlikely to

116    modify core linguistic capabilities in such a short time period. One explanation is that this

117    intervention could help PWA make more adaptive use of their current language system by

118    reducing maladaptive responses to language impairments related to extralinguistic factors such

119    as "linguistic anxiety" (Cahana-Amitay et al., 2011).

120         A key consideration when seeking to address system calibration is that some aspects of

121    language use may be more malleable and open to adaptive deployment than others. Speed-

122    accuracy tradeoffs have shown potential for malleability and evidence for adaptation deficits in

123    aphasia and are therefore worth pursuing from a rehabilitation perspective.

124

125    *System calibration, processing speed, and speed-accuracy tradeoffs in anomia*

126    Anomia is a word retrieval deficit measurable both in terms of accuracy and processing

127    speed (Moineau et al., 2005). Improving word retrieval accuracy has been the focus of most

128    anomia work to date (e.g., Best et al., 2013; Fridriksson et al., 2005), although speed has also

129    been considered in a number of instances (e.g., Neto & Santos, 2012; Prather et al., 1997).

130    However, the interactive relationship between speed and accuracy has not been adequately

131    considered in anomia treatment. In speed-accuracy tradeoffs, spending more time on a task tends

132    to increase accuracy, while spending less time lowers accuracy. Speed-accuracy tradeoffs are a

133    robust and widespread phenomenon in both human psychology (e.g., Wickelgren, 1977) and

134    beyond (Ceccarini et al., 2020). In humans, speed-accuracy tradeoffs appear to be partially under

135    volitional control: individuals are able to flexibly adjust speed vs. accuracy in the context of

136    shifting task instructions, feedback, or rewards that prioritize speed or accuracy (Campanella et

137    al., 2016; Starns & Ratcliff, 2010; Touron et al., 2007; Wagenmakers et al., 2008). Critically,

138    speed-accuracy tradeoffs are often nonlinear (Starns & Ratcliff, 2010), such that overly cautious

139    responses may be much slower but provide only marginal gains in accuracy, while overly

140    impulsive responses may be faster but result in much lower accuracy performance (Figure 2). In

141    previous response time modeling work, we have shown that speed-accuracy tradeoffs are present

142    in PWA during lexical decision and picture naming tasks. In Evans et al. (2019), we applied the

143    Diffusion Model (Ratcliff, 1978) to lexical decision data from 20 PWA, and found that 40%

144    demonstrated adaptation deficits in speed-accuracy tradeoffs, with impaired speed or accuracy

145    performance attributable to overly impulsive or overly cautious responses. In subsequent work

146    (Evans et al., 2020), we developed a novel multinomial ex-gaussian response time model of

147    picture naming in aphasia to estimate an "optimal response time cutoff," the point at which

148  additional processing time was unlikely to produce additional gains in accuracy. We fit this

149  model to picture naming data from PWA,

150  and found that for 8/10 participants, their

151  average response time (RT) for incorrect

152  responses exceeded their own optimal

153  RT cutoff. Together, these results

154  suggest that PWA do not always set

155  speed-accuracy tradeoffs to optimize task

156  performance in language-dependent

157  tasks.

158        If present, maladaptive speed-

159  accuracy tradeoffs likely have negative

160  consequences for everyday

161  communication and for treatment

162  outcomes. In everyday communication,

163  impulsive responses increase the chances



**Figure 2.** Schematic for speed-accuracy tradeoffs. Overly cautious responses slow response time (RT, red line) without further increasing accuracy (blue line). Overly impulsive responses improve RT, but with considerable lower levels of accuracy. Good system calibration (yellow vertical bar) balances these extremes and improves overall performance efficiency. Figure modified from Evans et al. (2019).

164  of making fast retrieval errors, self-corrections, and conversation repairs. On the other hand,

165  responses that are too cautious maximize accuracy in everyday communication at the cost of

166  timely transfer of information, slowed processing, and may make online communication

167  processes more susceptible to competition from internal or external distractions (e.g., PWA

168  forgetting their idea before they can finish sharing it).

169        In treatment, maladaptive speed-accuracy tradeoffs may have specific negative

170  consequences for dosage. Previous literature, particularly in the anomia context, has found that

171    dose-form typically consists of the number of successful/accurate retrieval events (Harvey et al.,

172    2020). Therefore, overly cautious responses would decrease dosage within a given treatment

173    time by decreasing the number of trials while only providing negligible additional gains in trial

174    accuracy. On the other hand, overly impulsive responses would increase error rates and

175    interference effects from error learning (Fillingham et al., 2006) without appreciably increasing

176    treatment dosage. As a result, PWA who display good system calibration may respond more

177    optimally during drill-based treatment – not too quickly, but willing to move on after their

178    chances of providing a correct response diminish – and as a result, maximize their treatment

179    dosage. Recent evidence within usual care from the VERSE trial suggests that treatment is most

180    successful when it is effortful while minimizing errors (but is not errorless), which provides

181    empirical evidence for this claim (Brogan et al., 2020). Thus, maladaptive speed-accuracy

182    tradeoffs in word retrieval likely have significant negative consequences for both everyday

183    communication effectiveness for PWA and negatively affect drill-based treatment dosage. In the

184    following section, we propose a general treatment framework intended to improve system

185    calibration as it relates to effort and speed-accuracy tradeoffs in the context of word retrieval.

186

187    *BEARS: A treatment framework to address strategic responses to core processing abilities*

188        Following the findings that PWA set maladaptive speed-accuracy tradeoffs and that

189    language performance in general may be affected by extralinguistic factors such as pressure,

190    frustration, or anxiety, we conceptualized an aphasia treatment framework called *BEARS*

191    (Balancing Effort, Accuracy, and Response Speed) intended to address these factors holistically.

192    The framework is intended to address PWA's underlying linguistic deficits through more

193    effective practice as well as increase strategic adaptation to underlying linguistic deficits through

194 increased awareness and training. Our claim underlying the BEARS framework is

195 straightforward: for an intervention to maximize its treatment dosage and its impact on everyday

196 communication, it should strive to strike a balance between processing effort, performance

197 accuracy, and response speed. Following the effortful retrieval practice literature (e.g., Middleton

198 et al., 2016) and findings that patient-generated responses are likely to compose key active

199 ingredients in treatment protocols (Evans et al., 2020; Gravier et al., 2018), responses within

200 treatment protocols should be effortful, but only up until the point where effort becomes

201 counterproductive.

202       Additionally, treatment should not only focus on constraining responses to this optimal

203 calibration point where effort, speed, and accuracy are balanced, but should provide explicit

204 metacognitive training and feedback to the person with aphasia, so that they are able to identify

205 this calibration point on their own without the need for input from a clinician. Instruction and

206 feedback can be implemented in terms of education on the relationship between these

207 components, metacognitive training including self-monitoring of frustration or tension which

208 reduce performance, and the successful use of strategies reduce these feelings which are

209 detrimental to successful performance. PWA should be encouraged to balance speed and

210 accuracy. Individuals who tend to make impulsive errors should be taught to slow down, while

211 individuals who tend to persevere and who are stuck for longer than is therapeutically ideal

212 should be taught to let unsuccessful attempts go and move on. Individuals who make both fast

213 and slow types of errors should be taught to notice and respond to both (see Appendix 1 for a

214 detailed description of the BEARS metacognitive training provided in the current study). Given

215 its nature, the BEARS treatment framework could be used to augment most evidence-based

216 restorative aphasia treatments targeting linguistic impairment, so that treatment not only

217　addresses the underlying impairment, but also a person with aphasia's strategic response to these

218　impairments. In the current study, we have applied the BEARS framework to one such treatment,

219　the semantic-feature verification (SFV; Kiran & Roberts, 2010).

220

221　*Semantic Feature Verification anomia treatment + BEARS*

222　　　　Semantically-oriented anomia treatments such as Semantic Feature Analysis (SFA;

223　Boyle, 2004; Boyle, 2010; Coelho et al., 2000) are among the most well-studied treatments for

224　naming impairment in aphasia. In SFA, the clinician shows the person with aphasia a pictured

225　object and elicits a naming attempt. In the traditional version of SFA (Boyle, 2010; Massaro &

226　Tompkins, 1994), the clinician then guides the PWA in verbally generating semantic features for

227　the target, using a chart specifying feature categories (Boyle & Coelho, 1995). Correct naming of

228　the target is elicited at the end of each trial. SFA has been modified in many ways since the

229　original papers, for example, to focus on verbs and actions (Wambaugh & Ferguson, 2007) or

230　implemented within the context of discourse (Peach & Reuter, 2010).  Kiran and Roberts (2010)

231　developed a variant of SFA in which repeated, guided practice centers around the verification of

232　semantic features for target words rather than their verbal production. As in generation-based

233　SFA, the verification-based variant (SFV) is hypothesized to improve retrieval of both treated

234　and untreated semantically-related words by strengthening the activation of related concepts in

235　the lexicon (e.g., Collins & Loftus, 1975). SFA has been found to improve treated words for

236　almost all participants and semantically related, untreated words for a large proportion of PWA

237　(Efstratiadou et al., 2018, Oh et al., 2016, Quique et al., 2018), including for the SFV variant

238　(e.g., Gilmore et al., 2020).

239    The current study employed a computer-based version of SFV as we were interested

240    evaluating the BEARS system calibration training in the context of an established anomia

241    treatment and the two-choice nature of SFV is well-suited for computer-based implementation

242    and feedback. The resulting BEARS + SFV treatment protocol included structured naming

243    practice and feature verification as well as education on speed-accuracy tradeoffs, metacognitive

244    training focused on the self-monitoring of effort, frustration, and timeliness of responses, and

245    computer-based performance feedback on the efficiency of both naming and feature verification

246    responses using a game-based points system (see Methods section). Thus, BEARS + SFA is not

247    only intended to strengthen the production of target words and underlying semantic networks

248    through more efficient drill-based practice, but also intended to improve participants' adaptive

249    system calibration, learning to make more adaptive use of their core language system during both

250    picture naming and feature-verification.

251

252    *Study purpose:*

253    The purpose of this study is to develop and pilot a BEARS-augmented anomia treatment

254    (BEARS + SFV) using a multiple baseline single-case series experimental design. Its goals are to

255    a) establish the feasibility of this approach, b) replicate previous SFV findings on performance

256    accuracy and determine whether BEARS + SFV improves naming and discourse production

257    efficiency, c) assess whether BEARS training improves how PWA respond in instances where

258    they cannot produce a target word,  and d) explore relationships between overall practice

259    efficiency and treatment outcomes. Positive findings will support further research developing

260    this intervention, which could establish comparative effectiveness of BEARS-augmented

261    compared to standard interventions.

262 *Study predictions:*

263     1. *BEARS + SFV will replicate previous SFA/SFV findings and improve naming accuracy*

264       *for both treated and semantically-related untreated words.*

265     2. *BEARS + SFV will increase naming efficiency.* By improving lexical access and system

266       calibration, BEARS + SFV will improve the efficient retrieval of trained and untrained

267       words, as measured in the number of correct words per minute.

268     3. *BEARS+SFV will improve discourse informativeness and efficiency.* This would indicate

269       that BEARS system calibration training generalizes beyond the single-word level where

270       it was trained.

271     4. *BEARS + SFV will improve system calibration for self-monitoring and error awareness.*

272       We predicted BEARS training would lead to a shift in the nature of how participants

273       responded on incorrect trials over time, producing a higher proportion of "pass"

274       responses and a corresponding reduction in overt errors (paraphasias) and timeout

275       nonresponses.

276     5. *Efficient practice performance during BEARS + SFV treatment will be positively*

277       *associated with good treatment outcomes.* While the current study cannot distinguish

278       correlation from causation, it is important to explore relationships between system

279       calibration, practice efficiency, and treatment outcomes to determine whether further

280       development of this work is warranted. We predicted that more efficient practice, would

281       be associated with larger treatment effect sizes.

282

283

284

285 **Methods**

286 *Participants*

287    Participants were recruited from the Western Pennsylvania Research Registry, the

288 Audiology and Speech Pathology Research Registry maintained by the VA Pittsburgh

289 Healthcare System (VAPHS), and local clinician referral. No participants enrolled in this study

290 received any concurrent speech-language treatment outside of the study-related sessions for the

291 duration of the study.

292    To be included in the study, participants were required to be at least 6 months post-onset

293 of stroke, have a diagnosis of aphasia (as defined by impairments in 2/8 subtest of the

294 Comprehensive Aphasia Test), be community-dwelling and at least wheel-chair ambulatory,

295 have spoken English as their primary language since childhood, and be age 18 or older.

296 Participants were also required to demonstrate less than or equal to 50% correct performance on

297 at least 80 treatment item probes during the pre-treatment study phase. Potential participants

298 were excluded if they had a history of neurodegenerative disease, active, unmanaged

299 psychopathology or alcohol/substance abuse, severe motor speech disorder (i.e., apraxia of

300 speech or dysarthria) or were participating in any other speech/language therapy during the time

301 of the study. Based on the complex multi-step nature of the treatment and our previous clinical

302 trial experience evaluating semantically-oriented anomia treatment (e.g., Evans et al., 2020), we

303 excluded participants who presented with very severe anomia, as measured by a CAT Naming

304 modality T-score of less than 40. Data collection took place at the VA Pittsburgh Healthcare

305 System with IRB approval (Study ID: Pro00002040).

306

307

15

308     *Assessment*

309         Participants were tested with standardized measures at study onset and again post-

310     treatment. They were assessed with the Comprehensive Aphasia Test (CAT; Swinburn et al.,

311     2004), the Philadelphia Naming Test (PNT; Roach et al., 1996), Cactus and Camel Test (CCT;

312     Bozeat et al., 2000), and selected subtests of the Psycholinguistic Assessments of Language

313     Processing in Aphasia (PALPA; Kay et al., 1996). Motor speech was assessed via the Duffy

314     protocol (Duffy, 2020) with diagnosis determined via consensus expert opinion between the first

315     and fourth authors who are certified SLPs.

316         Changes in monologue-based discourse informativeness and efficiency were evaluated

317     through the Nicholas and Brookshire protocol (Nicholas & Brookshire, 1993), which includes

318     two sets of discourse stimuli, each with two picture descriptions, one narrative, one procedural,

319     and one personal story. Discourse informativeness was measured by calculating the proportion of

320     correct information units (CIUs; effectively, words that are both accurate and relevant) to total

321     words. Efficiency was measured by the number of CIUs produced during the cumulative time

322     taken for each narrative task within (CIUs/minute). Calculation of CIUs, words, and time

323     followed the protocol described by Nicholas and Brookshire (1993). Informativeness and

324     efficiency is known to be relatively equivalent between sets and reasonably stable between

325     administrations (Nicholas & Brookshire, 1993). Sets were ordered pseudo-randomly. Samples

326     were scored by the treating, certified speech-language pathologist (fourth author) who is well-

327     trained in scoring CIUs but was not blinded to timepoint.

328

329

330

*Stimuli matching and selection*

332       Treatment and probe stimuli for study participants consisted of picturable nouns from two

333    freely available photographic databases (Brodeur et al., 2010; Brodeur et al., 2014; Moreno-

334    Martínez & Montoro, 2012). For each stimulus, we collected linguistic characteristics from

335    available corpora (Balota et al., 2007; Brysbaert et al., 2012) consisting of lexical frequency,

336    number of phonemes, and age of acquisition (Kuperman et al., 2012). Potential trained and

337    untrained items were then matched for production complexity based on an item complexity

338    algorithm from Fergadiotis et al. (2015). In this approach, item complexity was estimated using

339    the following equation: $B = -1.22 - .36(\log \text{word frequency}) + .21(\text{Age of Acquisition}) +$

340    $.15(\text{number of phonemes})$, which they reported to account for 63% of variance in naming

341    difficulty. This complexity score was used to create difficulty-matched triplets of items: a trained

342    item, an untrained related item from the same category, and an untrained unrelated item from a

343    different semantic category. Item triplets that had item complexity difference scores above 2

344    standard deviations were removed. The final set of stimuli consisted of 224 item triplets with

345    potential trained items across 15 semantic categories: body parts, building, clothing, decoration,

346    electronics, food, fruits and vegetables, furniture, kitchen utensils, mammals, nature, outdoor

347    activities, stationary, tools, and vehicles.

348       For each participant, treatment lists were generated on the basis of performance on a

349    confrontation picture-naming task. Pictures of the 224 potential treatment targets were presented

350    one at a time and participants were given 15 seconds to name each picture. Accuracy was judged

351    based on the 'first complete response' as per the scoring rules on the PNT. Per these rules, self-

352    corrections were not accepted if they already made a complete response as indicated by pausing

353    and/or prosody. If a participant indicated they did not know what the picture was, the item was

354    marked as incorrect but not selected as a treatment target. Each participant completed the naming

355    task on two separate occasions. Items that were named with less than or equal to 50% accuracy

356    across both administrations were included as potential treatment items. To qualify for treatment,

357    a category had to have at least 8 qualifying items. A total of 5 categories with 8 qualifying items

358    in each category were selected for treatment for a total of 40 treatment targets per participant.

359    Selection took into account participants personal interests and the quality of stimuli pictures.

360    Since each of the treatment targets had difficulty-matched related and unrelated generalization

361    items, administering all items would have created probe lists of 120 items. To reduce the

362    considerable testing burden, only one generalization item was randomly selected for each

363    treatment target, leading to a probe list of 80 items (40 treated words, 20 related untreated words,

364    and 20 unrelated untreated words).

365         For the semantic feature verification portion of the treatment task, eight semantic

366    "yes/no" feature questions for each of the 224 potential treatment targets were created by

367    undergraduate lab volunteers and were rated by 3 independent raters on a 1-10 scale (with "10"

368    being a good question and "1" being a poor question). Questions with an average score below

369    eight were re-written and rescored by three additional independent raters. For each treatment

370    target, four questions had a "yes" response and four had a "no" response. All questions were

371    audio-recorded and edited using Audacity software.

372

373    *Probe Administration*

374         Probes were administered for each participant in a multiple baseline across participants

375    design. Baseline probe performance was established via multiple probe assessments where each

376    participant was randomly assigned a number of baseline probe sessions (3, 4, or 5 sessions),

377   which helped control for the direct effects of probe exposure in the absence of treatment. Treated

378   and untreated items were assessed at the beginning of each treatment session prior to initiating

379   treatment, and within 1 week of finishing treatment. Items were also probed in a single follow-up

380   session approximately 3 weeks after the completion of treatment.

381        Probe administration during baseline, treatment, and follow-up included both

382   confrontation picture naming and a written lexical decision task, with item presentation

383   randomized within each task. For naming probes, participants were given 15 seconds to name

384   each picture and accuracy was judged based on the 'first complete response' as per stimuli

385   selection. The written lexical decision task was 320 trials in length, presenting all 80 probe

386   words and 80 matched pseudowords twice, and was always presented after naming probes each

387   session. Written lexical decision probes were collected for secondary response time modeling

388   analyses and are therefore not reported here.

389        Both probes and treatment software were programmed in PsychoPy software (Peirce et

390   al., 2019) and administered on a Dell XPS13 laptop using a USB microphone headset.

391   Assessment audio recordings were collected on a Surface Pro laptop using an external USB

392   microphone.

393

394   *Measuring naming response times*

395        Given the focus placed on speed-accuracy tradeoffs and efficiency in the current work,

396   trial response times during naming probes and treatment were collected via software voice key

397   and clinician button press (marked immediately after the first complete response had been

398   provided). Voice key responses were used to provide online computer-based points feedback

399   during treatment (see below), and therefore participants were trained to produce a single verbal

400  response and were reminded of these instructions each session. Voice key sensitivity was scored

401  online by the treating clinician, with apparent false or failed triggers noted by a key press. All

402  naming responses were also audio recorded for later review to establish rater reliability.

403      To assess voice key accuracy and reliability on naming probes, a trained independent

404  third rater hand-coded 10% of trials which the clinician had marked were triggered appropriately.

405  To do this, they viewed the recording waveforms in Audacity software and measured the

406  distance between stimuli onset and the first complete response by hand. Trials with large

407  disagreement between voice key and hand-coding were reviewed by the study team, and more

408  than 90% of these trials were due to ambiguity in PNT scoring rules for determining the first

409  correct response. Reliability between the voice key and hand-coded response times was good for

410  6 of the 9 participants ($r$'s ranging between .92 and .99), with poorer reliability for the remaining

411  3 participants ($r = .83$ for participant 1, $r = .24$ for participant 4, and $r = .69$ for participant 9).

412      However, after assessing the reliability of trials with 'good' voice key triggers above,  we

413  determined that naming trials with failed voice key triggers appeared to exclude data not at

414  random for some participants, as incorrect naming attempts were much more likely to be marked

415  as inaccurate voice key response times due to early partial production attempts. As a result, we

416  chose to use the clinician button-press measure of total trial time (from stimulus onset to

417  immediately after participants gave their first complete response) to inform measure of reward

418  rate used as a dependent variable in our efficiency analyses (see below).

419

420  *Treatment procedures*

421    Participants each received 25 hours of treatment administered by a licensed speech-

422    language pathologist (mostly by the 4[th] author[1]). Sessions were typically scheduled 3-4 days per

423    week with 2-3 hours of treatment time per day interspersed by breaks. Participants received

424    *BEARS+ SFV*, a hybrid clinician-computer treatment with 2 core components: computer-based

425    semantic feature verification treatment with points feedback, and BEARS meta-cognitive system

426    calibration training from a clinician. While being introduced to the treatment task, each

427    participant was educated on the speed-accuracy tradeoff and how to appropriately balance effort,

428    accuracy, and response speed. Participants were encouraged to find the balance and the "right

429    speed" for their processing ability in that moment. They were taught to become more aware of

430    instances when they are very unlikely to produce a correct response, and instructed to say "pass"

431    instead of producing an overt error or waiting until the response deadline had run out. This was

432    intended to reduce error learning and to increase the number of successful completed trials

433    during treatment.  Additional details regarding BEAR meta-cognitive strategy training and how

434    it was individualized for each participant are described in Appendix 1.

435    The treatment game was an implementation of SFV anomia treatment. Each treatment

436    trial consisted of 3 steps (Figure 3), with a naming attempt (step 1) followed by four feature

437    verification questions (step 2), followed by a second naming attempt (step 3).

438    In step 1, the target picture was presented, and the participant was asked to name it with a

439    single verbal response. Production accuracy, voice key success, and response time was judged

440    online by the clinician and input by key press. If the voice key trigger was successful,

441    participants then received immediate feedback consisting of their response accuracy and RT.

---

[1] The first author co-treated with the 4th for the first 5-6 sessions for the first 2 participants to ensure consistent application of the BEARS components. The first author was actively consulted and answered questions for the remainder of participants. The second author covered 1 session for participant 6 while the treating clinician was on vacation. All three are certified speech-language pathologists.

442        In step 2, the target picture was shown again along with an audio recording of the correct

443    response, followed by the auditory and written presentation of a semantic feature question.

444    Participants gave a yes/ no response to

445    each question via key press. Immediate

446    accuracy and response time feedback was

447    provided after they answered each

448    semantic feature question. Four of the

449    eight semantic feature questions were

450    randomly selected for presentation on each

451    treatment trial. After completing four

452    feature verification questions, Step 3 was

453    initiated.

454        In step 3, the participant was asked

455    to name the picture again "as quickly and

456    accurately as possible," and accuracy and

457    RT feedback were provided. Steps 1-3

458    were repeated for each target until all eight

459    category items had been practiced, which

460    completed a "round." At this point

461    cumulative point-based feedback for the

462    round was provided based on speed-



**Figure 3.** BEARS + SFV treatment schematic of computer-administered components.

463    accuracy performance, and then the process was repeated in a new round for a new category.

464    Every time all 40 items across the five categories were practiced, participants were told they had

465    completed a "level up" and point-based feedback was provided cumulative across all five rounds.

466    With- and between-category presentation order were determined randomly.

467         To support BEARS system calibration training, participants received both points

468    feedback and metacognitive strategy training throughout the intervention. Point feedback was

469    based on a modified reward rate algorithm (Bogacz et al., 2006) which rewarded efficiency (both

470    fast and accurate) performance, with correct responses earning 5 points and error responses

471    losing 1 point, divided by the total time spent on each treatment block. Points were awarded

472    separately for naming efficiency ("Coins") and feature verification efficiency ("Stars"). Points

473    were presented on the round and level up screens and the clinician reviewed the points with the

474    participant, including comparisons to previous blocks and treatment sessions to ensure

475    comprehension of the feedback system. Metacognitive strategy training throughout treatment

476    included ongoing education on the speed-accuracy tradeoff and participants' level of effort and

477    frustration.  Discussion included how speed can occasionally result in more errors, but how

478    slowing down does not always result in retrieving the target word. The clinician reviewed how

479    speech can be modulated and adjusted based on perceived difficulty/ accessibility of each target

480    and introduced the option to "pass" or "move on" once they recognize the feeling that the target

481    word is not accessible at the given time. Adaptive speed-accuracy tradeoffs were reinforced for

482    participants throughout treatment. A detailed description of this training and which participants

483    received which types of training are in Appendix 1 and Supplementary Materials S1.

484

485

486

487

488    *Analysis*

489         Data were analyzed in R statistical software, version 4.0.2 (R Core Team, 2020). Naming

490    probe data at all baseline and treatment timepoints were used for analysis of treatment response

491    for predictions 1, 2, and 4. For prediction 1 (*BEARS + SFV will replicate previous SFV findings*

492    *and improve naming accuracy for both treated and semantically-related untreated words*),

493    treatment outcomes were evaluated using the dependent variable of naming probe accuracy.  For

494    prediction 2 (*BEARS + SFV will increase naming efficiency*), treatment outcomes were evaluated

495    using the dependent variable of naming probe reward rate (i.e., the number of correct responses

496    per minute). For prediction 3 (*BEARS + SFV will improve discourse informativeness and*

497    *efficiency*), treatment outcomes were evaluated using the dependent variables of CIUs/minute

498    and proportion of CIUs. For prediction 4 (*BEARS + SFV will improve system calibration for*

499    *self-monitoring and error awareness*), self-monitoring ability was evaluated using the dependent

500    variable of naming probe "pass rate" (i.e., the proportion of inaccurate trials where participants

501    indicated they could not produce the target by saying "pass" instead of producing an overt error

502    or giving no response by the end of the 15-second response window). For prediction 5 (*Efficient*

503    *practice performance during BEARS + SFV treatment will be positively associated with good*

504    *treatment outcomes)*, treatment practice efficiency was evaluated using the dependent variable of

505    the total number of feedback points earned across treatment sessions ("coins" for naming

506    performance and "stars" for feature verification question performance), while dependent

507    variables for treatment outcomes were measured in terms of individual treatment effect sizes.

508         For predictions 1, 2, and 4, group-level performance was evaluated using Bayesian

509    generalized linear mixed-effect models using the R package BRMS (Bürkner, 2017) following

510    the interrupted time series approach described by Huitema and McKean (2000) and Moeyaert et

511   al. (2017). Bayesian implementations of generalized linear mixed effect models are largely

512   similar to their frequentist variant, but also permit estimation of the probability of a given effect

513   and individual effect sizes, as discussed below. The interrupted time series approach includes

514   fixed effects for baseline slope, level change, and slope change. Together, these fixed effects can

515   characterize the presence of a stable, rising, or declining baseline (i.e., baseline slope), any

516   immediate changes in performance at the onset of treatment (i.e., level change), and whether or

517   not the slope of treatment-related change exceeds that of the slope established during the baseline

518   phase (i.e., slope change). Therefore, positive level change and slope change fixed effects may

519   be found even in the presence of rising baseline slope, which provides evidence that changes in

520   probe performance over time are attributable to the treatment. Models were implemented

521   separately for each item condition (i.e., treated, related untreated, unrelated untreated). Details

522   regarding modeling fitting are reported in Appendix 2.

523        Individual effect sizes for naming accuracy were estimated for each participant by taking

524   the difference between the model's posterior prediction for each subject at the last treatment

525   probe and final baseline session, resulting in an estimate of the median number of words

526   improved and associated 90% credible interval. An equivalent approach was used to estimate the

527   individual improvements in naming reward rate. Group-level effect sizes were estimated by

528   calculating the difference between posterior samples at session 13 from session 4 at the mean of

529   the random effects for an average PWA, which accounts for performance during baseline.

530        A major benefit of this Bayesian mixed-effect approach is that a single model can

531   estimate effect sizes and group-level fixed effects, calculating 90% credible intervals and

532   posterior probabilities (i.e., the probability that the effect size or model parameter is greater than

533   zero) in each instance. Together, this a) provides an interpretable point estimate and range for

534  expected treatment effects, b) characterizes the degree of statistical robustness for effect sizes

535  based on posterior probabilities, and c) provides an appropriately conservative model check of

536  whether effect sizes are attributable to the treatment, which is done by comparing posterior

537  probabilities for the fixed effects of baseline slope, slope change, and level change.

538       For prediction 3, changes in discourse efficiency were calculated for each participant as a

539  measure of far generalization on the Nicholas and Brookshire protocol by calculating CIUs/min

540  and the proportion of CIUs . Changes in both proportion of CIUs and CIUs/minute were

541  analyzed via non-parametric bootstrap test for paired differences (Dwivedi et al., 2017) using the

542  R package *infer* (Bray et al., 2020). This approach is advantageous for small sample sizes as it

543  does not rely on underlying assumptions typical of parametric statistical tests and also

544  demonstrates better power than non-parametric tests in smaller sample sizes.

545       Prediction 5 was evaluated via Pearson correlations exploring relationships between

546  practice efficiency and treatment outcomes by condition (as calculated for predictions 1 through

547  4 above). Correlations between pass rates and treatment outcomes by condition were also

548  evaluated in the same correlation matrix as secondary analyses for prediction 4.

549

550  **Results**

551       Thirteen people with chronic aphasia were enrolled, and nine of the 13 met enrollment

552  criteria during initial assessment. At the time of their enrollment, participants ranged from 9

553  months to 44 years post-onset of a left-hemisphere stroke. Participants were between ages 55-73

554  years and were all native speakers of English. Apraxia of speech was either absent or very mild

555  and indistinguishable from phonological output deficits in the majority of participants, with the

556  exceptions of participant 1, who presented with moderate apraxia of speech and participant 2

557    who presented with probable mild apraxia of speech. Baseline testing on the Comprehensive

558    Aphasia Test and demographic information is presented in Tables 1 and 2, and additional

559    assessment results are reported in supplementary materials (S2).

560

561    **Table 1.** Participant Descriptive Characteristics

562

| Participant | Sex | Age | Etiology | MPO | Years Edu | Premorbid handedness | Race/ ethnicity | Hemiparesis |
|---|---|---|---|---|---|---|---|---|
| p1 | M | 63 | CVA | 259 | 14 | R | African American | R UE |
| p2 | M | 73 | CVA | 194 | 14 | R | Caucasian | None |
| p3 | M | 68 | CVA | 21 | 12 | L | Caucasian | None |
| p4 | M | 70 | CVA | 522 | 13 | R | Caucasian | R UE |
| p5 | M | 70 | CVA | 39 | 16 | R | Caucasian | None |
| p6 | F | 71 | CVA | 8 | 14+ years | R | African American | None |
| p7 | M | 70 | CVA | 9 | 18+ years | R | Caucasian | R UE and LE |
| p8 | M | 54 | CVA | 18 | 16 | R | African American | None |
| p9 | M | 72 | CVA | 58 | 14 | R | Caucasian | None |

563
564    *Note:* P= participant, M = Male, F = Female, CVA = Cerebrovascular accident, UE = Upper
565    Extremity, LE = Lower Extremity. P1 had concomitant moderate apraxia of speech and P2 had
566    probable mild apraxia of speech. All participants had a diagnosis of aphasia following left-
567    hemisphere CVA.
568

569

**Table 2.** Comprehensive Aphasia Test Performance (modality T score)

| Participant | Semantic Memory | Recognition Memory | Comp of Spoken Language | Comp of Written Language | Repetition | Naming | Reading (aloud) | Writing |
|---|---|---|---|---|---|---|---|---|
| p1 | 60 | 48 | 52 | 51 | 47 | 48 | 45 | 51 |
| p2 | 60 | 59 | 48 | 53 | 47 | 53 | 52 | 50 |
| p3 | 51 | 48 | 55 | 58 | 55 | 59 | 57 | 60 |
| p4 | 60 | 59 | 53 | 53 | 52 | 54 | 48 | 54 |
| p5 | 60 | 48 | 55 | 50 | 53 | 49 | 53 | 52 |
| p6 | 60 | 59 | 57 | 55 | 52 | 54 | 52 | 57 |
| P7 | 40 | 59 | 44 | 35 | 60 | 46 | 44 | 47 |
| P8 | 60 | 48 | 49 | 51 | 47 | 53 | 50 | 58 |
| P9 | 60 | 59 | 49 | 61 | 46 | 55 | 49 | 52 |

572

573      For predictions 1 and 2, aggregate naming accuracy and reward rate performance for each

574    participant by session are presented in Figure 4. For predictions 1, 2, and 4, model fixed effect

575    estimates and posterior probabilities are presented in Table 3 by dependent variable and item

576    condition, with full model results in Appendix 2.  For predictions 1 and 2, group and individual

577    effect sizes for naming accuracy and reward rate are presented in Figure 5.

578

579    *Results for prediction 1: BEARS + SFV will replicate previous SFV findings and improve*

580    *naming accuracy for both treated and semantically-related untreated words.*

581      Fixed effect beta-coefficients are presented as odds ratios. For treated items, credible

582    intervals excluded zero for baseline slope ($\beta$ = 0.20; 90% credible interval (CI) = 0.10, 0.29),

583    level change ($\beta$ = 0.56, 90% CI = 0.15, 0.98), and the quadratic term for slope change ($\beta$ = -0.02;

584    90% CI = -0.03, -0.01) but not slope change ($\beta$ = 0.03; 90% CI = -0.10, 0.15). For the untreated

585    conditions, only a positive trend of baseline slope ($\beta$ = 0.19, 90% CI = 0.08, 0.30) and downward

586    trend for slope change ($\beta$ = -0.17; 90% CI=-0.28, -0.05) for untreated items was evident. For

587     both untreated conditions, the quadratic term for slope change did not improve model fit and was

588     therefore not included.



**Figure 4.** Individual performance on naming probe accuracy and reward rate over time.

589     Calculation of overall group effect sizes revealed that an average PWA would be

590     expected to name an additional 13.7 treated words (90% CI: 7.3, 19.17) and demonstrate small

591     but meaningful generalization of 2.31 additional related, untreated words (90% CI: 0.31,4.36).

592     No improvements were seen for unrelated words (Effect size: 0.67, 90% CI: -1.37, 2.87). Group

593     effect sizes were consistent with individual effect size estimates (Figure 5). For treated items,

594     individual effect sizes ranged between two and 25 additional words named accurately, and the

595     90% credible intervals for excluded zero for all participants except for participant 7. The 90%

596     credible intervals excluded zero for four participants for untreated, related items, and excluded

597     zero for two participants for unrelated untreated items.

598     Table 3 provides a summary of treatment effect sizes by item condition along with fixed

599     effects and posterior probabilities to guide in their interpretation. Taken together, results showed

600     that participants improved on treated words (+13.7 additional words named accurately, posterior

601     probability > 0.99) and these gains were attributable to the treatment (posterior probability for

602     level change = 0.98) despite the presence of rising baselines (posterior probability > 0.99).

603     BEARS + SFV also produced improvements on related untreated words (+2.31 additional words

604     named accurately, posterior probability = 0.97), but these gains could not be clearly

605     distinguished from repeated probe exposure (level change and slope change posterior

606     probabilities ≤ .65). Finally, there were no improvements on unrelated untreated words (+0.67

607     additional words named accurately, posterior probability = 0.71).

608     These results provide clear evidence of a treatment effect for trained words and some

609     inconclusive evidence of response generalization to semantically-related untreated words, which

610     is broadly consistent with prediction 1.

611

612 **Table 3.** Summary of group treatment effect sizes, fixed effects, and posterior probabilities for
613 naming probe accuracy, reward rate, and "pass" rate.
614

| Model | | Treated words | | Related, Untreated words | | Unrelated, untreated words | |
|---|---|---|---|---|---|---|---|
| | | *Value* | *PP* | *Value* | *PP* | *Value* | *PP* |
| *Accuracy* | | | | | | | |
| | Effect Size | 13.70 | 1.00 | 2.31 | 0.97 | 0.67 | 0.71 |
| | Baseline Slope | 0.20 | 1.00 | 0.06 | 0.81 | 0.19 | 1.00 |
| | Level Change | 0.56 | 0.98 | 0.08 | 0.65 | -0.28 | 0.11 |
| | Slope Change | 0.03 | 0.63 | -0.01 | 0.44 | -0.17 | 0.01 |
| *Reward Rate* | | | | | | | |
| | Effect Size | 5.23 | 1.00 | 2.83 | 1.00 | 1.52 | 0.96 |
| | Baseline Slope | 0.22 | 1.00 | 0.16 | 1.00 | 0.27 | 1.00 |
| | level change | 0.49 | 0.99 | -0.03 | 0.43 | -0.32 | 0.03 |
| | slope change | -0.15 | 0.00 | -0.09 | 0.02 | -0.22 | 0.00 |
| *Pass Rate* | | | | | | | |
| | Effect Size | 0.24 | 0.98 | 0.24 | 0.99 | 0.18 | 0.99 |
| | Baseline Slope | -0.14 | 0.13 | 0.01 | 0.53 | -0.06 | 0.37 |
| | level change | 1.73 | 0.99 | 1.13 | 0.99 | 1.70 | 0.99 |
| | slope change | 0.55 | 1.00 | 0.14 | 0.80 | 0.51 | 0.98 |

615
616 *Note:* results are summarized from the Bayesian interrupted times series models completed for predictions 1 through
617 3. PP = The posterior probability (i.e., the probability that the effect size or model parameter is greater than zero).
618 Effect sizes are reported in additional words named correctly for accuracy, in additional words named per minute for
619 reward rate, and in increased proportion of "pass" vs. other inaccurate trial types for pass rate. Model fixed effects of
620 baseline slope, level change, and slope change are reported as odds ratios. Full model results are reported in
621 Appendix 2.
622

**Figure 5.** Group and Individual effect sizes for naming probe accuracy (top) and reward rate (bottom). Error bars reflect 90% Bayesian credible intervals.

623

624

625

626     *Results for prediction 2: BEARS + SFV will increase naming efficiency (reward rate).*

627         For treated items, there was a positive baseline slope ($\beta = 0.22$; 90% CI = 0.15, 0.29)

628     along with a positive level change ($\beta = 0.49$; 90% CI = 0.17, 0.83) and a negative slope change

629     ($\beta = -0.15$; 90% CI = -0.22, -0.09), suggesting that reward rate improved during baseline,

630     responded initially to treatment, and then slowed in its rate of improvement. For related and

631     unrelated untreated items, a rising baseline slope was evident (related: $\beta = 0.16$; 90% CI = 0.08,

632     0.23; unrelated: $\beta = 0.27$; 90% CI = 0.20, 0.33) but level change and slope change were not

633     meaningfully different from zero, indicating that the rate of improvement during treatment did

634     not exceed the slope established at baseline.

635         Calculation of overall group effect sizes revealed that the average participant's reward

636     rate improved by 5.23 additional words/minute for treated words (90% CI: 2.34, 10.99), and

637     demonstrated small gains of 2.83 additional words/minute (90% CI: 1.04, 6.25) for related

638     untreated words and 1.52 additional words/minute (90% CI: -0.06, 4.10) for unrelated untreated

639     words. Group effect sizes were consistent with individual effect size estimates (Figure 5). For

640     treated items, individual effect sizes ranged between 1.04 and 17.72 additional words named per

641     minute, and 90% credible intervals excluded zero for all nine participants. For semantically-

642     related untreated items, the 90% credible intervals excluded zero for all participants except for

643     participant 1. For unrelated untreated items, the 90% credible intervals excluded zero for five

644     participants.

645         As summarized in Table 3, results show that participants improved in naming efficiency

646     for treated words (+5.23 additional words/minute; posterior probability > 0.99) with gains

647     attributable to the treatment (posterior probability for level change = 0.99) despite the presence

648     of rising baselines (posterior probability > 0.99). Participants also demonstrated small, though

649    statistically robust, improvements in naming efficiency for related untreated words (2.83

650    additional words/minute; posterior probability > 0.99), and unrelated untreated words (1.52

651    additional words/minute; posterior probability = 0.96) but given rising baseline performance

652    (baseline slope posterior probabilities > 0.99) and lack of an increasing level or rate of

653    improvement during the treatment phase (level change and slope change posterior probabilities ≤

654    .43) these effects could not be clearly distinguished from repeated probe exposure.

655          The treatment-related gains in naming efficiency for treated items provides clear

656    evidence in support of prediction 2, while the inconclusive evidence for efficiency gains on

657    untreated items does not.

658

659    *Results for prediction 3: BEARS + SFV will improve discourse informativeness and efficiency.*

660          Performance on the Nicholas and Brookshire discourse elicitation task is reported in

661    Table 4. On average, participants improved their CIUs/minute by 4.07 and their proportion of

662    CIUs by 3.93% between entry and exit. However, these differences were not significant for

663    either CIUs/minute ($p = .18$) or proportion of CIUs ($p = .13$). Changes in discourse efficiency

664    (Table 4) were highly variable, with pre-post improvement by up to 21.29 CIUs/minute

665    (participant 7) and decreases as large as 12 CIUs/minute (participant 5[2]). These results do not

666    support prediction 3.

667

668

669

670

---

[2] Participant 5's discourse performance was highly tangential, which appeared to increase the overall session-to-session variability in his performance.

**Table 4. Pre- and Post-treatment Discourse Results (Nicholas and Brookshire protocol)**

| | CIUs/minute Pre-treatment | CIUs/minute Post-treatment | Total CIUs Pre-treatment | Total CIUs Post-treatment | % CIUs Pre-treatment | % CIUs Post-treatment | Total words Pre-treatment | Total words Post-treatment |
|---|---|---|---|---|---|---|---|---|
| p1 | 24.4 | 26.72 | 194 | 167 | 37.09 | 39.02 | 523 | 428 |
| p2 | 30.82 | 29.26 | 187 | 118 | 29.12 | 33.81 | 642 | 349 |
| p3 | 66.92 | 74.31 | 377 | 379 | 40.1 | 59.4 | 940 | 638 |
| p4 | 50.11 | 58.81 | 228 | 248 | 48.1 | 56.11 | 474 | 442 |
| p5 | 30 | 18 | 213 | 182 | 25.09 | 13 | 849 | 1399 |
| p6 | 20.06 | 25.73 | 237 | 193 | 53.26 | 53.17 | 445 | 363 |
| p7 | 38.2 | 59.49 | 198 | 235 | 52.94 | 63.17 | 374 | 372 |
| p8 | 21.73 | 30.45 | 151 | 101 | 28.76 | 34.12 | 525 | 296 |
| p9 | 42.22 | 38.29 | 254 | 224 | 47.65 | 45.71 | 533 | 490 |

*Results for prediction 4: BEARS + SFV will improve system calibration for self-monitoring and error awareness (pass rate).*

There was no effect of baseline slope for treated or untreated items (Appendix 2 and Table 3). A positive level change was present in all three conditions while a slope change was evident in the treated and unrelated conditions, suggesting that participants increased their proportion of "pass" responses compared to other error types after the onset of treatment and as treatment progressed. However, pass rates were highly variable and did not change equally for all participants (supplementary materials, S3).

As summarized in Table 3, results show that BEARS + SFV produced improvements in pass rate for treated words (mean rate increase of 0.24; posterior probability = 0.98) attributable to the treatment (posterior probability for level change = 0.99 and for slope change > 0.99). BEARS + SFV also produced similar changes in pass rate for related, untreated words (mean rate increase of 0.24; posterior probability = 0.99), and unrelated untreated words (mean rate increase 0.18; posterior probability = 0.99). Because of the stable baselines and high posterior

690    probabilities for slope and level change across conditions, we attribute changes in pass rate

691    directly to BEARS + SFV, indicating that the participants were able to increase their pass rate in

692    response to training. These results support prediction 4.

693            As a secondary analysis, we looked at correlations between pass rates, treatment effect

694    sizes, and aphasia severity (Figure 6). There was a particularly strong relationship between pass

695    rate and treated item effect sizes for naming accuracy ($r$ = -0.93), such that pass rates were

696    higher for participants with lower effect sizes. Correlations between pass rate and other naming

697    probe effect sizes were also negative ($r$s between -0.39 and -0.67), and the correlation between

698    pass rates and aphasia severity was weak ($r$ = -0.31). Visual inspection of pass rates by

699    participant over time (supplementary materials S3) shows that the three participants with the

700    lowest effect sizes for treated naming probe accuracy (participants 1, 5, and 7, Figure 5)

701    demonstrated the largest increases in pass rates as a result of treatment. These results are also

702    broadly consistent with prediction 4.

703

704    *Results for prediction 5: Efficient practice performance during BEARS + SFV treatment will be*

705    *positively associated with good treatment outcomes.*

706            Pearson correlations between practice efficiency (i.e., total feedback points earned with

707    "coins" for naming performance and "stars" for feature verification), naming probe treatment

708    effect sizes for accuracy and reward rate, and pre-post changes in discourse performance are

709    presented in Figure 6. Individual practice efficiency performance (coins and stars earned by

710    session) is presented in supplemental materials (S3). There was a strong positive correlation

711    between total coins and stars earned ($r$ = 0.94) and these measures therefore demonstrated

712    similar relationships to other variables. For correlations between practice efficiency and naming

713     probe effect sizes, there were a) strong positive correlations for reward rate on untrained related

714     items (coins, $r = 0.94$; stars, $r = 0.89$), b) strong positive correlations for reward rate on treated

715     items and on accuracy for related untreated items ($r$s between 0.67 and 0.73), and c) weak

716     positive correlations on accuracy for treated items (coins, $r = 0.31$; stars, $r = 0.25$). There were

717     weak positive correlations between practice efficiency and measures of discourse improvement

718     ($r$s between 0.22 and 0.33).

719         Aphasia severity has been found to predict treatment outcomes for semantically-oriented

720     anomia treatment (Quique et al., 2019) and could potentially play a role in how well individuals

721     were able to efficiently practice BEARS + SFV. Therefore, we examined correlations between

722     aphasia severity, practice efficiency, and treatment outcomes to explore whether practice

723     efficiency served as a proxy measure of aphasia severity in the current data (Figure 6). There

724     were moderate correlations between aphasia severity and practice efficiency (coins, $r = 0.37$;

725     stars $r = 0.42$), weak-to-moderate correlations between aphasia severity and naming probe effect

726     sizes ($r$s between 0.29 and 0.47), and essentially null correlations between aphasia severity and

727     discourse-related changes (proportion of CIUs, $r = 0.028$; CIUs per minute, $r = -0.14$).

728         Overall, these exploratory analyses demonstrated strong positive correlations between

729     practice efficiency and most measures of naming probe effect size, and only moderate

730     correlations between overall aphasia severity and practice efficiency, which is broadly consistent

731     with prediction 5.

**Figure 6.** Correlations and scatterplots between treatment outcomes, measures of treatment practice efficiency, and aphasia severity. "Coins" = total feedback points earned for naming practice efficiency, "Stars" = total feedback points earned for feature verification practice efficiency. "Acc ES: tx" = effect size for naming probe accuracy on treated items. "Acc ES: rel" = effect size for naming probe accuracy on semantically-related untreated items."RR ES: tx" = effect size for naming probe reward rate on treated items."RR ES: rel" = effect size for naming probe reward rate on semantically-related untreated items items. %CIU = proportion of CIUs. CAT mean = aphasia severity (CAT mean modality T score). r values > .67 are significant at uncorrected alpha = .05. r values > .836 are significant at Bonferroni-corrected alpha = .005.

732

733

734

**Discussion:**

The purpose of this study was to develop and pilot a BEARS-augmented anomia treatment (BEARS + SFV) that combined computer-based feedback and clinician-provided metacognitive system calibration training. We developed BEARS + SFV based on the rationale that training PWA to better balance their effort, accuracy and response speed during drill-based anomia treatment would allow them to make more efficient use of their *current* language system for the purposes of both language adaptation and restorative treatment. We hypothesized that this in turn could directly improve performance efficiency during probe tasks via a) better adaptation (e.g., PWA learning to inhibit avoidable paraphasias and/or move on more quickly from unproductive word retrieval attempts) and b) better restorative treatment outcomes for trained and related untrained words by allowing for by higher dosage, more successful effortful trials, and fewer produced errors during the course of treatment. Isolating the unique effects of BEARS-augmented compared to un-augmented anomia treatment was outside the scope of the current study as this would require a larger-scale comparative effectiveness study. BEARS+SFV showed good feasibility as a pilot treatment. The nine participants who met eligibility criteria successfully completed all study procedures and the full 25 hours of treatment, with zero attrition or loss to follow-up. The remaining study goals are addressed in turn.

*Does BEARS + SFV replicate previous SFV findings on performance accuracy and improve naming efficiency?*

Analyses for predictions 1 and 2 showed similar results for treated and untreated items on naming probes. Direct training effects were found for treated words, measured both in terms of accuracy (prediction 1) and reward rate (prediction 2). Rising baselines were noted in both

758    analyses, but additional treatment-related level increases in performance were noted in both

759    instances. Group effect sizes indicated that after 25 hours of BEARS + SFV training, a typical

760    individual with aphasia would be expected to accurately name an additional 13.7 out of 40

761    trained words and improve their efficiency by naming an additional 5.23 accurate trained words

762    per minute. At the individual level, eight out of nine participants improved on accuracy and all

763    nine improved in terms of reward rate. The accuracy effect size compares favorably to

764    previously reported anomia studies, which tend to train fewer words, often for equal or longer

765    periods of time (e.g., Snell et al., 2010). Together, results provide preliminary evidence that

766    participants improved in both accuracy and efficiency of treatment-related words as a result of

767    BEARS + SFV.

768        There was less robust evidence for treated-related gains for untrained words. Group effect

769    sizes showed clear improvement for semantically-related untreated words on naming probe

770    accuracy (consistent with prediction 1), and clear gains on naming probe reward rate for both

771    related and unrelated untreated words (consistent with prediction 2). However, rising baselines

772    were noted, and there was essentially no evidence of treatment-related slope or level change for

773    either the accuracy or reward rate analyses (posterior probabilities ≤ 0.65). In other words, while

774    participants improved their accuracy on related untreated words and their efficiency on both

775    related and unrelated untreated words over the course of the study, these gains cannot be clearly

776    distinguished from the positive effects of repeated exposure observed in the baseline phase[3].

---

[3] That being said, the calculation of fixed effects for treatment-related level and slope change rely on the modeling assumption that rising baselines would have continued linearly throughout the duration of the study even in the absence of treatment, which may be overly conservative. This highlights a general issue in single-subject designs when rising baselines are observed: it is not possible to determine whether a linear trend would have continued at the same rate (instead of leveling off) in the absence of treatment. This is problematic in regards to anomia treatment research in general, as repeated exposure to naming opportunities without feedback has been shown to improve naming abilities in at least some individuals (Creet et al., 2019). Given these considerations, we

777    These results broadly support prediction 1 ("*BEARS + SFV will replicate previous SFV*

778    *findings and improve naming accuracy for both treated and semantically-related untreated*

779    *words*"), and  are generally consistent with previous work looking at SFA and SFV (e.g., Kiran

780    & Roberts, 2010, Gilmore et al., 2020). We found robust direct training effects and weak

781    evidence of generalization to related untrained words. While rising baselines and methodological

782    issues (see Limitations) may have made generalization harder to capture, these results are also

783    consistent with previous claims that SFV may produce more modest generalization effects than

784    feature generation-based SFA (Boyle, 2010).

785    These results also partially supported prediction 2 ("*BEARS + SFV will increase naming*

786    *efficiency (reward rate)*"). There was clear evidence of treatment-related improvements on

787    naming efficiency observed across participants. Since reward rate is a measure of the total

788    number of correct responses per unit of time, some portion of efficiency gains for trained words

789    are attributable to improved word retrieval engendered by the SFV treatment component.

790    However, the BEARS system calibration component also appeared to play some role in these

791    gains, because even participant 7 who did not show clear gains in naming accuracy effect size

792    improved in reward rate effect size.

793    We predicted that if the BEARS treatment component improved naming efficiency on its

794    own via adaptive system calibration, participants would demonstrate efficiency gains on

795    *unrelated* untreated items, since the SFV treatment component was not predicted to improve

796    performance on these words. There were positive group-level gains in reward rate effect sizes for

797    unrelated untrained words (displayed by five of the nine participants at the individual level,

---

will interpret positive effect sizes in the absence of robust treatment-related slope and level change as weak positive
evidence (as opposed to null evidence) in the current study.

798 Figure 5), but as noted above, these gains could not be distinguished from rising baselines/

799 effects of repeated probe exposure based on the group-level fixed effects (Table 3). Overall,

800 these results provide clear evidence that BEARS + SFV improves naming efficiency for trained

801 words, but provides only weak evidence that BEARS training improves naming efficiency for

802 untrained words in some participants.

803

804 *Does BEARS+SFV improve discourse informativeness and efficiency?*

805 We hypothesized that improved system calibration induced BEARS training at the single-

806 word treatment during treatment could generalize to the level of connected speech and thereby

807 improve discourse informativeness and efficiency (prediction 3). However, this study did not

808 find any group-level changes in discourse performance (Table 4). While some individuals

809 improved in discourse informativeness and efficiency, others did not, and changes in these

810 patterns were not correlated with individual treatment response or aphasia severity (Figure 6).

811 These findings are consistent with the generally mixed findings regarding discourse-related

812 changes as a result of SFA (Rider et al., 2008; Silkes et al., 2020; Wallace & Kimelman, 2013).

813 Previous work has also shown test-retest instability on the Nicholas and Brookshire discourse

814 protocol (Cameron et al., 2010). This makes it difficult to detect treatment-related changes, and

815 suggests our discourse analyses were likely underpowered. In addition, BEARS+SVF computer

816 practice was completely focused on the single-word level, as was the great majority of the

817 complementary BEARS training. Therefore, the current null results suggest that BEARS training

818 at the single-word level is insufficient to induce system calibration improvements at the

819 discourse level (see future directions).

820

821 *Does BEARS + SFV improve system calibration for self-monitoring and error awareness?*

822    As part of training, participants were educated about speed-accuracy tradeoffs and how to

823    appropriately balance effort, accuracy, and response speed based on their own processing ability.

824    They were trained to provide a single naming response and to become more aware of instances

825    of anomia when they were very unlikely to produce a word correctly. In such instances, they

826    were instructed to say "pass" instead of producing an overt error or waiting until the allotted time

827    ran out, which was intended to increase overall efficiency and reduce the production of overt

828    errors. Therefore, we hypothesized that participants would demonstrate improved adaptation and

829    system calibration as by improving their proportion of "pass" responses relative to other error

830    types (paraphasias and timeout nonresponses). Results fully supported this prediction. As a

831    group, participants demonstrated improved pass rates attributable to the treatment. Correlation

832    analyses showed that pass rates were inversely proportional to treatment effect sizes for naming

833    probe accuracy, with higher pass rates for participants with lower effect sizes. This indicates

834    participants who demonstrated only small restorative treatment gains in response to the

835    restorative SFV component were still able to demonstrate strategic adaptation to their own poor

836    performance. Increased pass rates do not appear to have caused lower treatment responses,

837    because across participants, overt error and "pass" responses took longer on average than correct

838    responses. This means that choosing to say "pass" was more likely to replace an overt error than

839    a correct response, and therefore did not reflect giving up too early and losing opportunities for

840    correct retrieval. Overall, we interpret these findings as evidence of improved system calibration

841    as a result of the BEARS component of this intervention.

842        Participant 7 provides a helpful illustrative example of these effects: although he was the

843    only participant who did not improve on probe accuracy (effect size 90% credible intervals all

844    included zero, Figure 5 top panel), he demonstrated small but robust gains in naming efficiency

845    across item categories (Figure 5, bottom panel). Given his steadily increasing pass rates

846    (supplementary materials, Figure S3a), his gains in naming probe efficiency are likely

847    attributable to improved system calibration as opposed to restored naming ability.

848    *Is practice efficiency positively associated with treatment outcomes?*

849         One of the primary rationales for developing BEARS + SFV was a consideration of how

850    overly impulsive and overly cautious speed-accuracy tradeoffs could negatively affect treatment

851    outcomes in drill-based restorative treatment. Overly impulsive responses may increase the

852    number of avoidable errors and therefore may reduce treatment outcomes via error learning

853    (Fillingham et al., 2006), while overly cautious responses are slow and may reduce treatment

854    outcomes via reduced overall dosage. Based on this premise, the computer-based feedback in

855    BEARS + SFV used an algorithm which awarded points based on maximizing the number of

856    correct responses while minimizing the number of errors over time. This feedback was designed

857    to help participants better balance speed and accuracy during practice in a way that maximized

858    both effortful and errorless learning principles (Schuchard & Middleton, 2018). We hypothesized

859    that more efficient practice, as measured by a great number of total "Coin" and "Star" feedback

860    points earned, would be correlated with larger treatment effect sizes (prediction 5). This

861    prediction was largely confirmed.

862         While our group-level analyses for prediction 1 only found weak support for

863    generalization to related untrained items (i.e., small effect sizes without clear changes in

864    treatment-related slope or level change), we actually found strong correlations between practice

865    efficiency and generalization to related untrained naming probe accuracy using a case series

866    correlational approach (Rapp, 2011). In other words, participants who produced more correct

867    naming and feature verification responses and fewer errors over 25 hours of treatment also

868    demonstrated better response generalization to semantically-related untrained items. In contrast,

869    there were only weak correlations between practice efficiency and trained naming probe

870    accuracy. We attribute this to the fact that best-possible accuracy performance for trained items

871    was observed approximately half-way through the treatment for most participants (Figure 4).

872    This suggests that dosage in the current study was sufficient to produce participant-specific

873    ceiling effects for trained item accuracy (regardless of individual practice efficiency), but that the

874    degree of response generalization specifically predicted by the SFV treatment component was

875    affected by overall levels of practice efficiency.

876          In addition, we found strong positive correlations between practice efficiency and

877    individual effects sizes for naming efficiency on both trained and related untrained words. It is

878    not particularly surprising that more efficient practice during treatment corresponded to more

879    efficient probe performance, but it does support the importance of considering not only accuracy

880    but also efficiency in anomia treatment.

881          We would also note that these exploratory findings related to practice efficiency do not

882    appear to be merely a result of aphasia severity. Correlations between aphasia severity, practice

883    efficiency and effect sizes suggest that a relatively small amount of variance was shared between

884    practice efficiency and aphasia severity ($r^2 = 0.15$), while higher proportions of variance were

885    shared between practice efficiency and related untrained naming probe treatment response ($r^2$s

886    ranging between 0.45 and 0.88). Individual predictors of this treatment response should be

887    further explored.  Overall, these results provide support for targeting practice efficiency during

888    drill-based restorative anomia treatment. However, these results are preliminary and cannot

889    distinguish correlation from causation. Future work should confirm these exploratory findings

890    and evaluate the relative contributions of efficiency-focused practice in relation to learning

891    theory.

892    *Preliminary guidance for treatment candidacy based on individual case results.*

893         Most of the participants who showed good treatment responses were those who improved

894    in their practice efficiency during the course of treatment, and clear increases in practice

895    efficiency were apparent by the fourth treatment session for these participants (Supplementary

896    Figure S3b). This suggests that early treatment performance may be predictive of overall

897    treatment response (Simic et al., 2020).

898         In contrast, the three participants with the lowest treatment responses in naming accuracy

899    (participants 1, 5, and 7) also had the most severe anomia per PNT scores and participant 7 had

900    the low semantic control performance of all participants as measured by the CCT

901    (Supplementary language testing: S2), and Participant 1 was the only individual in the sample

902    who presented with clear (more than suspected/very mild) apraxia of speech. Participants 5 and 7

903    were also the only two participants who consistently presented with overly conservative response

904    patterns during treatment (see Appendix 1 and Supplementary materials S1). These

905    characteristics may be negatively prognostic in terms of overall treatment response. However,

906    some degree of nuance is required in considering this characterization of "non-responders".

907         Participants 1, 5 and 7 all demonstrated minimal naming gains naming probe accuracy,

908    but all three improved in proportion of "pass" compared to other error responses, demonstrated

909    modest improvements in naming efficiency, and participants 5 and 7 also demonstrated gains in

910    functional communication per anecdotal family report (Supplementary materials S1). Together,

911    these patterns suggest that future research should consider the distinct effects of restorative vs.

912    compensatory BEARS training, as PWA who are poor restorative treatment candidates may still

913    be good candidates for improved adaptation and system calibration.

914

915    *Study limitations:*

916         There were limitations in this study regarding probe design, stimuli selection, stimuli

917    scoring, participant selection, and dosing of treatment components that should be addressed in

918    future work. Despite being matched for general production difficulty using an algorithm based

919    on word frequency, age of acquisition, and word length in phonemes (Fergadiotis et al., 2015),

920    untreated items had higher baseline performance than treated items, which may have negatively

921    affected our ability to detect treatment-related generalization. We administered potential

922    treatment items twice for selection (and selected items that were ≤ 50% accurate), but relied on

923    difficulty matching to select untreated items in order to decrease testing burden (i.e., 224 instead

924    of 672 words, administered twice). In retrospect, we realize that our approach assumed that if a

925    participant could not name an easier potential treatment word, then they would be equally

926    unlikely to name other words of the same difficulty level. However, 37% of the variance in

927    naming accuracy was unexplained in Fergadiotis et al.'s (2015) prediction model, and our pre-

928    selection assessment of treated but not untreated words in essence created a filter that allowed for

929    "regression to the mean" on untrained words. As a result, when we selected easier words for

930    treatment, participants were more likely to be able to name their matched pairs correctly.

931    Predicted treatment generalization was still observed in terms of group effect sizes, but could not

932    be distinguished from effects of repeated probe exposure; these effects may be partially

933    attributable to this design limitation. As an additional limitation, naming and discourse probes

934    were scored by the fourth author (who administered most of the assessment and treatment) and

935    therefore were not blinded for time-point.

936        As discussed above, simple probe exposure can produce rising baseline effects (Creet et

937    al., 2019). However, an additional potential limitation in this study was that PWA saw the

938    written word-form for all treated and untreated items in a lexical decision task that was presented

939    at each probe timepoint, which was administered for the purposes of secondary analyses not

940    reported here. The lexical decision task was always provided after naming probes to minimize

941    potential bias, and no feedback was provided in either task, but positive effects of this additional

942    exposure cannot be ruled out. However, if simple repeated probe exposure in picture plus written

943    form without feedback is indistinguishable from treatment-related generalization resulting from

944    25 hours of intensive SFV treatment, then the clinical efficiency of SFV and similar

945    semantically-orient anomia treatments may require reevaluation. Given the pilot nature of this

946    study, we chose not to recruit PWA with very severe anomia given their low treatment response

947    to SFA (Quique et al., 2019). However, our participants who did not improve in terms of

948    accuracy still improved in efficiency and pass rate. Therefore, the adaptive component of

949    BEARS may be effective for individuals with more severe anomia and this should be examined

950    in future work.

951        As noted, BEARS + SFV consists of two core components: computer-based semantic

952    feature verification naming treatment and BEARS meta-cognitive system calibration training.

953    While the dosage of meta-cognitive component was flexibly adjusted in the current study based

954    on perceived participant need, the relative dosage of each component needs to be more carefully

955    characterized in future comparative effectiveness work applying clearly established treatment

956     fidelity procedures. Clinicians seeking to apply this approach should also weigh the relative

957     benefits of each component for a given patient, based on their ongoing performance.

958         In aligning the current pilot with previous SFV work, several additional methodological

959     differences should be noted when interpreting these findings. We trained a larger number of

960     words than is typical (40 items split between 5 categories). While previous approaches have

961     trained individual categories sequentially, we trained all five categories within each session, with

962     practice blocked by category. However, these differences did not appear to be detrimental based

963     on favorable treated effect sizes for trained words. In addition, the efficiency focus of the

964     BEARS treatment component could have decreased depth of processing for the semantic feature

965     questions. But if so, faster processing of questions would have decreased generalization effects

966     to semantically-related untrained items, and the opposite pattern was found. Finally, treatments

967     that rely on SFV also generally manipulate category typicality (e.g., Gilmore et al., 2020), which

968     was not specifically addressed here. Given these differences, the current study findings reflect a

969     modest extension of SFV-based anomia treatment.

970

971     *Future directions:*

972         In the current study, we relied on measures of efficiency (e.g., reward rate for naming

973     probes) as proxy measures for system calibration, but future work could apply response time

974     modeling techniques (Evans et al., 2019; Evans et al., 2020) to better assess changes in overly

975     impulsive and overly conservative speed-accuracy tradeoffs and provide individualized

976     computer-based feedback.

977         BEARS + SFV did not produce reliable changes in discourse informativeness or

978     efficiency at the group level, suggesting that the word-level focus of this treatment was too far

979 removed as a practice context to promote generalization to connected speech (Stokes & Baer,

980 1977; Thompson, 1989). However, there was a great deal of individual variability and several

981 additional findings that suggest at least some changes may have occurred beyond the single-word

982 level. These including anecdotal family reports of improved functional communication

983 (participants 3, 5, and 7) and at least one instance of negative training transfer to conversation

984 consisting of the overgeneralization of "pass" responses (participant 4, Supplementary materials

985 S3). Therefore, it is worth exploring whether BEARS system calibration training might better

986 improve communication efficiency at the discourse level (Whitney & Goldstein, 1989) if used to

987 augment treatments which target discourse performance directly (e.g., BEARS + Attention

988 Reading and Constrained Summarization, Rogalski & Edmonds, 2008).

989 When providing BEARS meta-cognitive training and feedback, we paid a great deal of

990 attention to participants' body language and visible muscle tension as proxies for their frustration

991 and overall effort. Electromyography or other appropriate biosignals may be sensitive to these

992 observations, which could allow for the development of biofeedback-based BEARS training.

993 In motivating BEARS + SFV, we drew a distinction between the PWA's core linguistic

994 impairments and how adaptively they make use of their current system to optimize performance.

995 It is likely that the ability to flexibly adjust to linguistic impairments and task constraints relies

996 on domain-general cognitive abilities, which have been increasingly implicated as concomitant

997 deficits in PWA (e.g., Gilmore et al., 2019; Murray, 2012). They also may depend on person-

998 level factors such as linguistic anxiety (Cahana-Amitay et al., 2011). Therefore, cognitive and

999 person-level predictors of treatment response should be explored in future work adequately

1000 powered to examine such effects.

1001    In addition, we used the conceptualization of adaptation vs. impairment symptoms from

1002    Adaptation Theory (Kolk & Heeschen, 1990) to motivate this study's focus on adaptive system

1003    calibration. Original support for Adaptation Theory reported differences in self-awareness of

1004    grammatical output for people with fluent vs. nonfluent aphasia, with nonfluent individuals with

1005    Broca's aphasia argued to be more aware of the grammaticality of their output than individuals

1006    with Wernicke's aphasia, and therefore more able to strategically adapt their spoken output using

1007    telegraphic speech (Kolk & Heeschen, 1992). While research based on aphasia classification is

1008    outside the current scope of this pilot study (and inconsistent with our choice to characterize the

1009    language profiles of our participants using the CAT), future work could explore the relationship

1010    between aphasia syndromes, error awareness, and response to BEARS system calibration

1011    training. This would be especially important if BEARS was used to augment discourse-level

1012    treatments, which would be more dependent upon fluency considerations.

1013

1014    *Conclusion:*

1015    The purpose of this study was to develop and pilot a BEARS-augmented anomia

1016    treatment (BEARS + SFV) that combined computer-based feedback and metacognitive system

1017    calibration training. BEARS + SFV showed good feasibility as a pilot treatment. Results

1018    provided strong evidence for direct training effects on naming probe accuracy and some weaker

1019    evidence of generalization to semantically-related untrained words, consistent with previous

1020    semantically-oriented anomia SFA/ SFV research (Kiran & Roberts, 2010, Quique et al., 2019,

1021    Boyle, 2010). Results provided strong evidence for direct training effects on naming probe

1022    efficiency, and some weaker evidence that the BEARS treatment component improved naming

1023    efficiency for untrained words. There were no group-level improvements in measures of

discourse performance, but participants did demonstrate improved system calibration based on their ability to shift the nature of their responses on inaccurate treatment trials, with an increasing proportion of "pass" responses compared to paraphasia or timeout nonresponses. In addition, computer-based feedback and BEARS training was designed to promote practice efficiency, and practice efficiency during treatment was positively correlated with treatment outcomes. Follow-up work will be necessary to replicate these effects and distinguish correlation from causation, but these findings are consistent with the claim that improving practice efficiency in SFV anomia treatment leads to greater treatment generalization and improved naming efficiency.

Overall, this study establishes the feasibility of BEARS + SFV and provides preliminary evidence that it improves naming efficiency, especially for trained words, and response adaptation for inaccurate trials. Therefore, future work should examine BEARS-augmented compared to standard aphasia interventions in well-powered comparative effectiveness research designed to characterize specific contributions of BEARS training on restorative and compensatory treatment outcomes.

On a final note, the current study also highlights the importance of considering processing speed in addition to accuracy in anomia treatment. People with very mild aphasia still report frustration over language efficiency (Cavanaugh & Haley, 2019), and need to improve in efficiency even when at ceiling for accuracy (Neto & Santos, 2012). On the other end of the proficiency spectrum, people with severe aphasia and non-responders may still have the potential to improve language efficiency, as suggested by the current findings. Inefficient communication is a major source of frustration, which makes it an important treatment target in its own right.

1045       While the current study provides only preliminary support for BEARS, applying elements

1046 of this system calibration treatment framework may still be of interest to clinicians, and a

1047 summary of BEARS training and education materials have been made available in Appendix 1.

1048

1060

**References**

Balota, D. A., Yap, M. J., Hutchison, K. A., Cortese, M. J., Kessler, B., Loftis, B., Neely, J. H.,
Nelson, D. L., Simpson, G. B., & Treiman, R. (2007). The English lexicon project.
*Behavior Research Methods, 39*(3), 445-459.

Best, W., Greenwood, A., Grassly, J., Herbert, R., Hickin, J., & Howard, D. (2013). Aphasia
rehabilitation: Does generalisation from anomia therapy occur and is it predictable? A
case series study. *Cortex, 49*(9), 2345-2357.
https://doi.org/https://doi.org/10.1016/j.cortex.2013.01.005

Bogacz, R., Brown, E., Moehlis, J., Holmes, P., & Cohen, J. D. (2006). The physics of optimal
decision making: a formal analysis of models of performance in two-alternative forced-
choice tasks. *Psychological Review, 113*(4), 700-765.

Boyle, M. (2004). Semantic Feature Analysis Treatment for Anomia in Two Fluent Aphasia
Syndromes. *American Journal of Speech-Language Pathology, 13*(3), 236-249.
http://dx.doi.org/10.1044/1058-0360(2004/025)

Boyle, M. (2010). Semantic feature analysis treatment for aphasic word retrieval impairments:
what's in a name? *Top Stroke Rehabil, 17*(6), 411-422. https://doi.org/10.1310/tsr1706-
411

Boyle, M., & Coelho, C. A. (1995). Application of semantic feature analysis as a treatment for
aphasic dysnomia. *American Journal of Speech-Language Pathology, 4*(4), 94-98.

Bozeat, S., Lambon Ralph, M. A., Patterson, K., Garrard, P., & Hodges, J. R. (2000). Non-verbal
semantic impairment in semantic dementia. *Neuropsychologia, 38*(9), 1207-1215.
https://doi.org/10.1016/s0028-3932(00)00034-8

1083 Bray, A., Ismay, C., Chasnovski, E., Baumer, B., Cetinkaya-Rundel, M., Couch, S., Laderas, T.,

1084          Solomon, N., Hardin, J., Kim, A. Y., Fultz, N., Friedman, D., Cotton, R., & Fannin, B.

1085          (2020). *infer Package.* In (Version 0.5.3) https://cran.r-

1086          project.org/web/packages/infer/infer.pdf

1087 Brodeur, M. B., Dionne-Dostie, E., Montreuil, T., & Lepage, M. (2010). The Bank of

1088          Standardized Stimuli (BOSS), a New Set of 480 Normative Photos of Objects to Be Used

1089          as Visual Stimuli in Cognitive Research. *PLOS ONE, 5*(5), e10773.

1090          https://doi.org/10.1371/journal.pone.0010773

1091 Brodeur, M. B., Guérard, K., & Bouras, M. (2014). Bank of Standardized Stimuli (BOSS) Phase

1092          II: 930 New Normative Photos. *PLOS ONE, 9*(9), e106953.

1093          https://doi.org/10.1371/journal.pone.0106953

1094 Brogan, E., Godecke, E., & Ciccone, N. (2020). Behind the therapy door: what is "usual care"

1095          aphasia therapy in acute stroke management? *Aphasiology*, 1-23.

1096          https://doi.org/10.1080/02687038.2020.1759268

1097 Brysbaert, M., New, B., & Keuleers, E. (2012). Adding part-of-speech information to the

1098          SUBTLEX-US word frequencies. *Behavior Research Methods, 44*(4), 991-997.

1099          https://doi.org/10.3758/s13428-012-0190-4

1100 Bürkner, P. C. (2017). Advanced Bayesian multilevel modeling with the R package brms. *arXiv*

1101          *preprint arXiv:1705.11123*.

1102 Cahana-Amitay, D., Albert, M. L., Pyun, S. B., Westwood, A., Jenkins, T., Wolford, S., &

1103          Finley, M. (2011). Language as a Stressor in Aphasia. *Aphasiology, 25*(2), 593-614.

1104          https://doi.org/10.1080/02687038.2010.541469

1105  Cameron, R. M., Wambaugh, J. L., & Mauszycki, S. C. (2010). Individual variability on

1106      discourse measures over repeated sampling times in persons with aphasia. *Aphasiology,*

1107      *24*(6-8), 671-684. https://doi.org/10.1080/02687030903443813

1108  Campanella, F., Skrap, M., & Vallesi, A. (2016). Speed-accuracy strategy regulations in

1109      prefrontal tumor patients. *Neuropsychologia, 82*, 1-10.

1110      https://doi.org/10.1016/j.neuropsychologia.2016.01.008

1111  Carragher, M., Conroy, P., Sage, K., & Wilkinson, R. (2012). Can impairment-focused therapy

1112      change the everyday conversations of people with aphasia? A review of the literature and

1113      future directions. *Aphasiology, 26*(7), 895-916.

1114      https://doi.org/10.1080/02687038.2012.676164

1115  Cavanaugh, R., & Haley, K. L. (2020). Subjective communication difficulties in very mild

1116      aphasia. *American Journal of Speech-Language Pathology*, *29*(1S), 437–448.

1117      https://doi.org/10.1044/2019_AJSLP-CAC48-18-0222

1118  Ceccarini, F., Guerra, S., Peressotti, A., Peressotti, F., Bulgheroni, M., Baccinelli, W., Bonato,

1119      B., & Castiello, U. (2020). Speed–accuracy trade-off in plants. *Psychonomic Bulletin &*

1120      *Review*. https://doi.org/10.3758/s13423-020-01753-4

1121  Coelho, C. A., McHugh, R. E., & Boyle, M. (2000). Semantic feature analysis as a treatment for

1122      aphasic dysnomia: A replication. *Aphasiology, 14*(2), 133-142.

1123  Collins, A. M., & Loftus, E. F. (1975). A spreading-activation theory of semantic processing.

1124      *Psychological Review, 82*(6), 407.

1125  Conroy, P., Sotiropoulou Drosopoulou, C., Humphreys, G. F., Halai, A. D., & Lambon Ralph,

1126      M. A. J. B. (2018). Time for a quick word? The striking benefits of training speed and

1127      accuracy of word retrieval in post-stroke aphasia. *141*(6), 1815-1827.

1128    Creet, E., Morris, J., Howard, D., & Nickels, L. (2019). Name it again! investigating the effects

1129        of repeated naming attempts in aphasia. *Aphasiology, 33*(10), 1202-1226.

1130        https://doi.org/10.1080/02687038.2019.1622352

1131    Dietz, A., Thiessen, A., Griffith, J., Peterson, A., Sawyer, E., & McKelvey, M. (2013). The

1132        renegotiation of social roles in chronic aphasia: Finding a voice through AAC.

1133        *Aphasiology, 27*(3), 309-325. https://doi.org/10.1080/02687038.2012.725241

1134    Dwivedi, A. K., Mallawaarachchi, I., & Alvarado, L. A. (2017). Analysis of small sample size

1135        studies using nonparametric bootstrap test with pooled resampling method. *Statistics in*

1136        *Medicine, 36*(14), 2187-2205. https://doi.org/10.1002/sim.7263

1137    Efstratiadou, E. A., Papathanasiou, I., Holland, R., Archonti, A., & Hilari, K. (2018). A

1138        Systematic Review of Semantic Feature Analysis Therapy Studies for Aphasia. *J Speech*

1139        *Lang Hear Res, 61*(5), 1261-1278. https://doi.org/10.1044/2018_jslhr-l-16-0330

1140    Evans, W. S., Hula, W. D., Quique, Y., & Starns, J. J. (2020). How much time do people with

1141        aphasia need to respond during picture naming? estimating optimal response time cutoffs

1142        using a multinomial ex-gaussian approach. *Journal of Speech, Language, and Hearing*

1143        *Research, 63*(2), 599-614. https://doi.org/doi:10.1044/2019_JSLHR-19-00255

1144    Evans, W. S., Hula, W. D., & Starns, J. J. (2019). Speed–Accuracy Trade-Offs and Adaptation

1145        Deficits in Aphasia: Finding the "Sweet Spot" Between Overly Cautious and Incautious

1146        Responding. *American Journal of Speech-Language Pathology, 28*(1S), 259-277.

1147        https://doi.org/10.1044/2018_AJSLP-17-0156

1148    Fergadiotis, G., Kellough, S., & Hula, W. D. (2015). Item response theory modeling of the

1149        Philadelphia Naming Test. *Journal of Speech, Language, and Hearing Research, 58*(3),

1150        865-877.

1151    Fillingham, J. K., Sage, K., & Lambon Ralph, M. A. (2006). The treatment of anomia using

1152    errorless learning. *Neuropsychological Rehabilitation*, *16*(2), 129-154. https://doi-

1153    org.pitt.idm.oclc.org/10.1080/09602010443000254

1154    Fridriksson, J., Holland, A. L., Beeson, P., & Morrow, L. (2005). Spaced retrieval treatment of

1155    anomia. *Aphasiology, 19*(2), 99-109. https://doi.org/10.1080/02687030444000660

1156    Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013).

1157    *Bayesian data analysis*. CRC press.

1158    Gilmore, N., Meier, E. L., Johnson, J. P., & Kiran, S. (2019). Non-linguistic cognitive factors

1159    predict treatment-induced recovery in chronic post-stroke aphasia. *Archives of Physical*

1160    *Medicine and Rehabilitation, 100*(1), 1251-1258.

1161    https://doi.org/10.1016/j.apmr.2018.12.024

1162    Gilmore, N., Meier, E. L., Johnson, J. P., & Kiran, S. (2020). Typicality-based semantic

1163    treatment for anomia results in multiple levels of generalisation. *Neuropsychological*

1164    *Rehabilitation, 30*(5), 802-828. https://doi.org/10.1080/09602011.2018.1499533

1165    Goodglass, H. (1980). Disorders of naming following brain injury. *American Scientist, 68*(6),

1166    647-655.

1167    Gravier, M. L., Dickey, M. W., Hula, W. D., Evans, W. S., Owens, R. L., Winans-Mitrik, R. L.,

1168    & Doyle, P. J. (2018). What Matters in Semantic Feature Analysis: Practice-Related

1169    Predictors of Treatment Response in Aphasia. *American Journal of Speech-Language*

1170    *Pathology, 27*(1S), 438-453.

1171    Hartig, F. (2020). *DHARMa: Residual diagnostics for hierarchical (multi-level/mixed)*

1172    *regression models* In

1173     Harvey, S. R., Carragher, M., Dickey, M. W., Pierce, J. E., & Rose, M. L. (2020). Treatment

1174         dose in post-stroke aphasia: A systematic scoping review. *Neuropsychol Rehabil*, 1-32.

1175         https://doi.org/10.1080/09602011.2020.1786412

1176     Heeschen, C., & Schegloff, E. A. (1999). Agrammatism, adaptation theory, conversation

1177         analysis: On the role of so-called telegraphic style in talk-in-interaction. *Aphasiology,*

1178         *13*(4-5), 365-405.

1179     Huitema, B. E., & McKean, J. W. (2000). Design Specification Issues in Time-Series

1180         Intervention Models. *Educational and Psychological Measurement, 60*(1), 38-58.

1181         https://doi.org/10.1177/00131640021970358

1182     Hunt, P., Soto, G., Maier, J., Müller, E., & Goetz, L. (2002). Collaborative teaming to support

1183         students with augmentative and alternative communication needs in general education

1184         classrooms. *Augmentative and Alternative Communication, 18*(1), 20-35.

1185         https://doi.org/10.1080/aac.18.1.20.35

1186     Kay, J., Lesser, R., & Coltheart, M. (1996). Psycholinguistic assessments of language processing

1187         in aphasia (PALPA): An introduction. *Aphasiology, 10*(2), 159-180.

1188     Kiran, S., & Roberts, P. M. (2010). Semantic feature analysis treatment in Spanish–English and

1189         French–English bilingual aphasia. *Aphasiology, 24*(2), 231-261.

1190     Kolk, H., & Heeschen, C. (1990). Adaptation symptoms and impairment symptoms in Broca's

1191         aphasia. *Aphasiology, 4*(3), 221-231.

1192     Kolk, H., & Heeschen, C. (1992). Agrammatism, paragrammatism and the management of

1193         language. *Language and Cognitive Processes*, *7*(2), 89–129.

1194         http://doi.org/10.1080/01690969208409381

1195    Kuperman, V., Stadthagen-Gonzalez, H., & Brysbaert, M. (2012). Age-of-acquisition ratings for

1196        30,000 English words. *Behavior Research Methods 44*(4), 978-990.

1197        https://doi.org/10.3758/s13428-012-0210-4

1198    Marshall, R. S., Laures-Gore, J., & Love, K. (2018). Brief mindfulness meditation group training

1199        in aphasia: exploring attention, language and psychophysiological outcomes.

1200        *International Journal of Language & Communication Disorders, 53*(1), 40-54.

1201        https://doi.org/10.1111/1460-6984.12325

1202    Massaro, M., & Tompkins, C. A. (1994). Feature analysis for treatment of communication

1203        disorders in traumatically brain-injured patients: An efficacy study. *Clinical aphasiology,*

1204        *22*, 245-256.

1205    Middleton, E. L., Schwartz, M. F., Rawson, K. A., Traut, H., & Verkuilen, J. (2016). Towards a

1206        Theory of Learning for Naming Rehabilitation: Retrieval Practice and Spacing Effects. *J*

1207        *Speech Lang Hear Res, 59*(5), 1111-1122. https://doi.org/10.1044/2016_jslhr-l-15-0303

1208    Moeyaert, M., Ugille, M., Natasha Beretvas, S., Ferron, J., Bunuan, R., & Van den Noortgate,

1209        W. (2017). Methods for dealing with multiple outcomes in meta-analysis: a comparison

1210        between averaging effect sizes, robust variance estimation and multilevel meta-analysis.

1211        *International Journal of Social Research Methodology, 20*(6), 559-572.

1212    Moineau, S., Dronkers, N. F., & Bates, E. (2005). Exploring the processing continuum of single-

1213        word comprehension in aphasia. *J Speech Lang Hear Res, 48*(4), 884-896.

1214        https://doi.org/10.1044/1092-4388(2005/061)

1215    Moreno-Martínez, F. J., & Montoro, P. R. (2012). An Ecological Alternative to Snodgrass &

1216        Vanderwart: 360 High Quality Colour Images with Norms for Seven Psycholinguistic

1217        Variables. *PLOS ONE, 7*(5), e37527. https://doi.org/10.1371/journal.pone.0037527

1218     Murray, L. L. (2012). Attention and Other Cognitive Deficits in Aphasia: Presence and Relation

1219         to Language and Communication Measures. *American Journal of Speech-Language*

1220         *Pathology, 21*(2), S51-S64. https://doi.org/10.1044/1058-0360(2012/11-0067)

1221     Neto, B., & Santos, M. E. (2012). Language after aphasia: Only a matter of speed processing?

1222         *Aphasiology, 26*(11), 1352-1361. https://doi.org/10.1080/02687038.2012.672023

1223     Nicholas, L. E., & Brookshire, R. H. (1993). A system for quantifying the informativeness and

1224         efficiency of the connected speech of adults with aphasia. *Journal of Speech, Language,*

1225         *and Hearing Research, 36*(2), 338-350.

1226     Oh, S. J., Eom, B., Park, C., & Sung, J. E. (2016). Treatment Efficacy of Semantic Feature

1227         Analyses for Persons with Aphasia: Evidence from Meta-Analyses. *Commun Sci Disord,*

1228         *21*(2), 310-323. https://doi.org/https://doi.org/10.12963/csd.16312

1229     Peach, R. K., & Reuter, K. A. (2010). A discourse-based approach to semantic feature analysis

1230         for the treatment of aphasic word retrieval failures. *Aphasiology, 24*(9), 971-990.

1231     Peirce, J., Gray, J. R., Simpson, S., MacAskill, M., Höchenberger, R., Sogo, H., Kastman, E., &

1232         Lindeløv, J. K. (2019). PsychoPy2: Experiments in behavior made easy. *Behavior*

1233         *Research Methods, 51*(1), 195-203. https://doi.org/10.3758/s13428-018-01193-y

1234     Prather, P., Zurif, E., Love, T., & Brownell, H. (1997). Speed of lexical activation in nonfluent

1235         Broca's aphasia and fluent Wernicke's aphasia. *Brain and Language, 59*(3), 391-411.

1236     Quique, Y., Evans, W. S., & Dickey, M. W. (2018). Acquisition and Generalization Responses

1237         in Aphasia Naming Treatment: A Meta-Analysis of Semantic Feature Analysis

1238         Outcomes. *American Journal of Speech-Language Pathology, 28*(1S), 230-246.

1239         https://doi.org/10.1044/2018_AJSLP-17-0155

1240    R Core Team. (2020). *R: A Language and Environment for Statistical Computing.* In (Version

1241        3.6.3) R Foundation for Statistical Computing. https://www.R-project.org

1242    Rapp, B. (2011). Case series in cognitive neuropsychology: Promise, perils, and proper

1243        perspective. *Cognitive Neuropsychology*, *28*(7), 435–444.

1244        http://doi.org/10.1080/02643294.2012.697453

1245    Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review, 85*(2), 59-108.

1246    Rider, J. D., Wright, H. H., Marshall, R. C., & Page, J. L. (2008). Using semantic feature

1247        analysis to improve contextual discourse in adults with aphasia. *American Journal of*

1248        *Speech-Language Pathology, 17*(2), 161-172.

1249    Roach, A., Schwartz, M. F., Martin, N., Grewal, R. S., & Brecher, A. (1996). The Philadelphia

1250        naming test: Scoring and rationale. *Clinical aphasiology, 24*, 121-133.

1251    Rogalski, Y., & Edmonds, L. A. (2008). Attentive Reading and Constrained Summarisation

1252        (ARCS) treatment in primary progressive aphasia: A case study. *Aphasiology, 22*(7-8),

1253        763-775. https://doi.org/10.1080/02687030701803796

1254    Schuchard, J., & Middleton, E. L. (2018). The roles of retrieval practice versus errorless learning

1255        in strengthening lexical access in aphasia. *Journal of Speech, Language, and Hearing*

1256        *Research*, *61*(7), 1700-1717.  https://doi:10.1044/2018_JSLHR-L-17-0352

1257    Silkes, J. P., Fergadiotis, G., Graue, K., & Kendall, D. L. (2020). Effects of Phonomotor Therapy

1258        and Semantic Feature Analysis on Discourse Production. *American Journal of Speech-*

1259        *Language Pathology*, 1-14. https://doi.org/https://doi.org/10.1044/2020_AJSLP-19-

1260        00111

1261    Simic, T., Chambers, C., Bitan, T., Stewart, S., Goldberg, D., Laird, L., Leonard, C., Rochon, E.

1262        (2020). Mechanisms underlying anomia treatment outcomes. *Journal of Communication*

1263        *Disorders*, *88*, 106048. http://doi.org/https://doi.org/10.1016/j.jcomdis.2020.106048

1264    Simmons-Mackie, N. (2018). *Aphasia in North America*. https://www.aphasiaaccess.org/state-of-

1265        aphasia

1266    Snell, C., Sage, K., & Lambon Ralph, M. A. (2010). How many words should we provide in

1267        anomia therapy? A meta-analysis and a case series study. *Aphasiology, 24*(9), 1064-1094.

1268        https://doi.org/10.1080/02687030903372632

1269    Starns, J., & Ratcliff, R. (2010). The effects of aging on the speed–accuracy compromise:

1270        Boundary optimality in the diffusion model. *Psychology and Aging, 25*(2), 377-390.

1271    Stokes, T. F., & Baer, D. M. (1977, Summer). An implicit technology of generalization. *Journal*

1272        *of applied behavior analysis, 10*(2), 349-367. https://doi.org/10.1901/jaba.1977.10-349

1273    Swinburn, K., Porter, G., & Howard, D. (2004). *Comprehensive Aphasia Test*. Psychology Press.

1274    Thompson, C. K. (1989). Generalization research in aphasia: A review of the literature. *Clinical*

1275        *aphasiology, 18*, 195-222.

1276    Touron, D. R., Swaim, E. T., & Hertzog, C. (2007). Moderation of Older Adults' Retrieval

1277        Reluctance Through Task Instructions and Monetary Incentives. *The Journals of*

1278        *Gerontology: Series B, 62*(3), P149-P155. https://doi.org/10.1093/geronb/62.3.P149

1279    Wagenmakers, E. J., Ratcliff, R., Gomez, P., & McKoon, G. (2008). A Diffusion Model Account

1280        of Criterion Shifts in the Lexical Decision Task. *J Mem Lang, 58*(1), 140-159.

1281        https://doi.org/10.1016/j.jml.2007.04.006

1282    Wallace, S. E., & Kimelman, M. D. (2013). Generalization of word retrieval following semantic

1283        feature treatment. *Neurorehabilitation, 32*(4), 899-913.

1284     Wambaugh, J. L., & Ferguson, M. (2007). Application of semantic feature analysis to retrieval of

1285          action names in aphasia. *Journal of Rehabilitation Research & Development, 44*(3).

1286     Whitney, J. L., & Goldstein, H. (1989). Using self-monitoring to reduce disfluencies in speakers

1287          with mild aphasia. *J Speech Hear Disord, 54*(4), 576-586.

1288          https://doi.org/10.1044/jshd.5404.576

1289     Wickelgren, W. A. (1977). Speed-accuracy tradeoff and information processing dynamics. *Acta*

1290          *Psychologica, 41*(1), 67-85. https://doi.org/10.1016/0001-6918(77)90012-9

1291

*Description of appendices:*

Appendix 1. Detailed description of clinician-provided BEARS metacognitive training

Appendix 2. Bayesian model fitting procedures and full model results


*Description of supplementary materials:*

**S1.** by-participant detailed description of individualized BEARS training.


**S2.** additional language testing results


**S3a.** Pass rate (proportion of "pass" responses relative to paraphasia and nonresponse error types) by session. P1, P5, and P7 were the three participants who responded least to the treatment in terms of naming probe accuracy.


**S3b.** Practice efficiency during treatment over time, as measured by cumulative feedback points by session.

**Appendix 1: Detailed description of clinician-provided BEARS metacognitive training**

**1. General education regarding balancing effort, accuracy, and response speed.** To begin, all participants were educated on the speed-accuracy tradeoff at the beginning of treatment. A figure similar to Figure 2 (but abstracted with the specific numbers removed) was used to supplement the discussion of speed and accuracy. Education, explaining the tradeoff between speed and accuracy and the eventual plateau of accuracy even with increasing processing time, was provided. Participants were introduced to the possibility of adjusting word-finding speed based on perceived difficulty/accessibility of each target. Participants were also taught to say "pass" or "move on" once they decided that the target word was not accessible in the moment. Additionally, they were instructed to only respond with one word in a single attempt. Cues to provide only one attempt were provided throughout treatment. Proper use of speed-accuracy tradeoff was explicitly reinforced for participants throughout treatment. Each participant received an individualized meta-cognitive portion of the treatment, in which both speed and effort were addressed. Clinician judgement was used to determine which meta-cognitive strategies in each category were most useful for an individual's treatment.

**2. Self-monitoring tension and frustration and learning when to relax.** The clinician subjectively monitored participants' muscle tension, body language, vocal quality, and prosody as proxy measures of retrieval effort, especially when noted in conjunction with unsuccessful word retrieval. If a participant appeared to be using too much effort in an attempt to access the target word, the clinician would draw the participant's attention to the feeling of effort and tension and recommended strategies to reduce the tension and frustration. The clinician would encourage the participant to notice when they were starting to get tense and frustrated, and to take a breath and/or deliberately relax tension in their neck and shoulders.

**3. Analogies to understand anomia and speed-accuracy tradeoffs.**
• *Example analogy 1*. Imagine a junk drawer in your garage that is filled to the brim with tools. If you want to get the screwdriver out of the drawer and just grab for something without taking your time, you will likely select a different tool like a hammer or a wrench. Alternatively, if you take just a little bit more time, you could easily select the screwdriver. Similarly, if you move too quickly to select a target word, it is more likely that you will select the wrong word or sounds. You need to take your time so that you are sure you have "the right tool".
• *Example analogy 2*. Tension and processing interference were likened to creating radio static in the head. Participants were told to "adjust the radio dial" to make the single clearer by relaxing and allowing for more efficient access.
• *Example analogy 3*. Building a case vs. making a snap judgement. The clinician highlighted the differences of building a case and making an immediate decision when you're in the moment. Naming a picture would be an immediate decision. If long delays were observed, the clinician would remind the participant that it is not about building a case, but rather making that quick decision in the moment.

**4. Not waiting too long.** If a participant was waiting more than 15 seconds, the clinician would provide education that increased wait time beyond a certain point will likely not result in successfully access of the target word. Additionally, the clinician would highlight the relationship between increased wait times and increased tension and frustration. The participant would be encouraged to "move on" as soon as they felt that they were having trouble accessing the target word. The clinician would lead the participant through discussion and reflection on what it feels like for them when they cannot access the target word in that moment. This feeling

would be contrasted with scenarios where the participant was able to access the target word within 1-10 seconds so they could learn when to "move on."

**5. Modulating speed based on retrieval difficulty.** If a participant began to overgeneralize speed (always slowing down or speeding up regardless of accuracy performance), the clinician would encourage them to modulate their speed based on perceived difficulty. For example, they may be encouraged to give themselves additional time to access a polysyllabic word such as "sledgehammer" as compared to a monosyllabic word such as "bit." Additionally, in these instances, they would be encouraged to slow down their articulation rate and give themselves the time to say the word accurately.

**6. Identifying the *right* amount of effort.** Some participants were engaged in a discussion about using the "*right* amount of effort." "Too much effort" can result in visible tension and participants may feel as though they are trying to "push out" a word. Alternatively, participants were encouraged to be sure they were engaged with enough effort so as not to "operate on cruise control." After a large number of trials, participants would occasionally repeat the last word they heard in the semantic feature verification question or the previous target. They were encouraged to monitor the amount of cognitive resources they were devoting to the task.

**7. Working *with* your aphasia.** Education was provided regarding the concept of "working *with* your aphasia" instead of against it. Discussion surrounded the idea that the participant will likely still experience word finding difficulty following treatment. However, responding adaptively will help avoid exacerbating those difficulties. They were encouraged to use strategies to manage their frustration, effort, and timing to increase the likelihood of accessing the target word.

**8. Distinguishing initiation speed vs. articulation rate.** Participants were encouraged to name pictures quickly, which was occasionally led to increased articulation rate as well. Since increasing articulation rate often resulted in additional sound errors, participants were instructed to give themselves time to say the word accurately once they had accessed it and to not overly speed up their articulation rate.

**9. Noting differences between conversation and the treatment task.** Some participants questioned the benefits of the training to "pass" or "move on" or were noted to be using this approach to some extent during conversation. Here, education focused on the benefits of identifying when word retrieval was going to be unsuccessful during both treatment and in conversation. During anomia treatment, participants should "move on" to keep practicing, while during conversation, this could be the same point to shift towards alternative communication strategies (e.g., writing the word instead of saying it).

**10. Wait until you have one word.** Encouraging speed was often met with semantic paraphasias after a number of successive trials. When this occurred, participants were encouraged to slow down and only produce one word to name the target. This helped to reduce the number of immediate self-corrections.

**Table A1. Summary of the BEARS metacognitive training strategies used with each participant.**

| | p1 | p2 | p3 | p4 | p5 | p6 | p7 | p8 | p9 |
|---|---|---|---|---|---|---|---|---|---|
| 1. General education regarding BEARS | x | x | x | x | x | x | x | x | x |
| 2. Self-monitoring tension/frustration and learning to relax | x | x | x | x | | x | x | x | x |
| 3. Analogies for anomia and speed-accuracy tradeoffs | x | | x | x | | x | x | | |
| 4. Not waiting too long | | | | | x | | x | | |
| 5. Modulating speed based on retrieval difficulty | | | x | x | | x | | x | x |
| 6. Identifying the *right* amount of effort | x | | | | | | x | x | x |
| 7. Working *with* your aphasia | x | x | x | | | x | x | x | x |
| 8. Distinguishing initiation speed vs. articulation rate. | x | x | | | | x | | x | x |
| 9. Differences between conversation and treatment. | | | | x | | | x | | x |
| 10. Wait until you have one word | | | | | x | x | | x | x |

**Appendix 2: Model Fitting Procedures and Full Model Results**

For predictions 1, 2, and 3, group-level performance was evaluated using Bayesian generalized linear mixed-effect models using the R package BRMS (Bürkner, 2017) following the interrupted time series approach described by Huitema and McKean (2000) and Moeyaert et al. (2017). Models were implemented separately for each item condition (i.e., treated, related untreated, unrelated untreated). The proportion of correct naming attempts was modeled using a binomial distribution via a logistic link function. Reward rate and "pass" rate were modeled using a lognormal distribution with a gaussian link function.

For each dependent variable (i.e., accuracy, reward rate, and "pass" rate) and item condition, a maximal model with random intercepts by participant and correlated random slopes for each predictor variable was initially fit. Each model was then reduced iteratively on the basis of posterior predictive checks, leave-one-out cross validation, and model convergence (in this case eliminating the number of divergent transitions). Because visualization of the aggregate data indicated that response to treatment may be nonlinear, a quadratic term for baseline session and slope change were also evaluated and only maintained in a given model if they significantly improved model fit. Differences between item conditions were determined by evaluating whether 90% credible intervals overlapped between conditions.

As reward rate is a combination of a count variable (the number of correct responses) and response time, the model was fit in an iterative fashion with multiple probability distributions, including the lognormal, truncated normal, and gamma distributions, with the lognormal distribution showing the best fit. Model fit for each of these probability distributions was compared via posterior predictive checks and leave-one-out cross validation.

For naming accuracy, models were evaluated for overdispersion via simulation using the R package DHARMa (Hartig, 2020); overdispersion was not present ($p > .05$). 4000 iterations were run for each of four independent Hamiltonian Markov Chain Monte Carlo; the initial 2000 chains were discarded and not included in the estimation of each parameter. Models were run with weakly informative priors: normal distributions with a mean of 0 and standard deviation of 10 for beta coefficients and a half-cauchy distribution with a mean of 0 and standard deviation of 5. Models were assessed for convergence using the split-half potential scale reduction factor (Gelman et al., 2013) and the effective sample size. In all the models the estimated split-half potential scale reduction factor values were less than 1.01, and the number of effective sample sizes exceeded 1000 for all parameters. Posterior predictive checks confirmed the models adequately fit the data.

**Table A1.** Naming Probe Accuracy Bayesian General Mixed-Effect Model Results

| Model | Parameter | Mean | Std. Error | Lower CI | Upper CI |
|---|---|---|---|---|---|
| | *Population level effects* | | | | |
| | Intercept | -1.89 | 0.26 | -2.32 | -1.48 |
| | Baseline Slope | 0.2 | 0.06 | 0.10 | 0.29 |
| | Level Change | 0.56 | 0.26 | 0.15 | 0.98 |
| | Slope Change | 0.03 | 0.08 | -0.10 | 0.15 |
| Treated | Slope Change$^2$ | -0.02 | 0.01 | -0.03 | -0.01 |
| | *Group level effects* | | | | |
| | sd: Intercept | 0.56 | 0.24 | 0.27 | 0.98 |
| | sd: Baseline Slope | 0.08 | 0.04 | 0.02 | 0.15 |
| | sd: Level Change | 0.51 | 0.26 | 0.16 | 0.97 |
| | sd: Slope Change | 0.05 | 0.04 | 0.00 | 0.13 |
| | | | | | |
| | *Population level effects* | | | | |
| | Intercept | -0.84 | 0.34 | -1.40 | -0.30 |
| | Baseline Slope | 0.06 | 0.07 | -0.05 | 0.16 |
| | Level Change | 0.08 | 0.22 | -0.28 | 0.45 |
| Related | Slope Change | -0.01 | 0.07 | -0.12 | 0.10 |
| | *Group level effects* | | | | |
| | sd: Intercept | 0.84 | 0.29 | 0.48 | 1.38 |
| | sd: Baseline Slope | 0.04 | 0.03 | 0.00 | 0.09 |
| | sd: Level Change | 0.17 | 0.14 | 0.01 | 0.45 |
| | sd: Slope Change | 0.04 | 0.03 | 0.00 | 0.10 |
| | | | | | |
| | *Population level effects* | | | | |
| | Intercept | -1.34 | 0.33 | -1.89 | -0.82 |
| | Baseline Slope | 0.19 | 0.07 | 0.08 | 0.30 |
| | Level Change | -0.28 | 0.23 | -0.64 | 0.10 |
| Unrelated | Slope Change | -0.17 | 0.07 | -0.28 | -0.05 |
| | *Group level effects* | | | | |
| | sd: Intercept | 0.79 | 0.29 | 0.44 | 1.33 |
| | sd: Baseline Slope | 0.05 | 0.03 | 0.01 | 0.10 |
| | sd: Level Change | 0.21 | 0.17 | 0.02 | 0.54 |
| | sd: Slope Change | 0.04 | 0.03 | 0.00 | 0.10 |

Note: Beta-coefficients are presented as log odds. CI refers to the 90% credible interval.

**Table A2.** Naming Probe Reward Rate Bayesian General Mixed-Effect Model Results

| Model | Parameter | Mean | Std. Error | Lower CI | Upper CI |
|---|---|---|---|---|---|
| | *Population level effects* | | | | |
| | Intercept | -0.17 | 0.27 | -0.60 | 0.26 |
| | Baseline Slope | 0.22 | 0.04 | 0.15 | 0.29 |
| | Level Change | 0.49 | 0.21 | 0.17 | 0.83 |
| Treated | Slope Change | -0.15 | 0.04 | -0.22 | -0.09 |
| | *Group level effects* | | | | |
| | sd: Intercept | 0.72 | 0.25 | 0.43 | 1.17 |
| | sd: Baseline Slope | 0.04 | 0.03 | 0.01 | 0.09 |
| | sd: Level Change | 0.48 | 0.21 | 0.20 | 0.85 |
| | sd: Slope Change | 0.03 | 0.03 | 0.00 | 0.08 |
| | | | | | |
| | *Population level effects* | | | | |
| | Intercept | 0.53 | 0.36 | -0.05 | 1.10 |
| | Baseline Slope | 0.16 | 0.04 | 0.08 | 0.23 |
| | Level Change | -0.03 | 0.15 | -0.27 | 0.22 |
| Related | Slope Change | -0.09 | 0.05 | -0.17 | -0.02 |
| | *Group level effects* | | | | |
| | sd: Intercept | 0.95 | 0.32 | 0.59 | 1.53 |
| | sd: Baseline Slope | 0.03 | 0.02 | 0.00 | 0.07 |
| | sd: Level Change | 0.14 | 0.11 | 0.01 | 0.35 |
| | sd: Slope Change | 0.03 | 0.02 | 0.00 | 0.07 |
| | | | | | |
| | *Population level effects* | | | | |
| | Intercept | 0.21 | 0.32 | -0.30 | 0.72 |
| | Baseline Slope | 0.27 | 0.04 | 0.20 | 0.33 |
| | Level Change | -0.32 | 0.17 | -0.60 | -0.04 |
| Unrelated | Slope Change | -0.22 | 0.04 | -0.29 | -0.15 |
| | *Group level effects* | | | | |
| | sd: Intercept | 0.86 | 0.28 | 0.52 | 1.40 |
| | sd: Baseline Slope | 0.03 | 0.02 | 0.00 | 0.06 |
| | sd: Level Change | 0.32 | 0.16 | 0.08 | 0.61 |
| | sd: Slope Change | 0.03 | 0.02 | 0.00 | 0.07 |

**Table A3.** Naming Probe "Pass" Rate Bayesian General Mixed-Effect Model Results

| Model | Parameter | Mean | Std. Error | Lower CI | Upper CI |
|---|---|---|---|---|---|
| Treated | *Population level effects* | | | | |
| | Intercept | -3.01 | 0.62 | -4.06 | -2.04 |
| | Baseline Slope | -0.14 | 0.13 | -0.37 | 0.07 |
| | Level Change | 1.73 | 0.78 | 0.46 | 3.01 |
| | Slope Change | 0.55 | 0.18 | 0.25 | 0.85 |
| | Slope Change$^2$ | -0.03 | 0.02 | -0.06 | -0.01 |
| | *Group level effects* | | | | |
| | sd: Intercept | 1.38 | 0.53 | 0.72 | 2.35 |
| | sd: Baseline Slope | 0.08 | 0.07 | 0.01 | 0.22 |
| | sd: Level Change | 1.85 | 0.63 | 1.05 | 2.99 |
| | sd: Slope Change | 0.20 | 0.13 | 0.03 | 0.44 |
| | sd: Slope Change$^2$ | 0.03 | 0.02 | 0.01 | 0.06 |
| Related | *Population level effects* | | | | |
| | Intercept | -2.36 | 0.50 | -3.20 | -1.58 |
| | Baseline Slope | 0.01 | 0.13 | -0.20 | 0.22 |
| | Level Change | 1.13 | 0.53 | 0.28 | 2.02 |
| | Slope Change | 0.14 | 0.17 | -0.14 | 0.41 |
| | Slope Change$^2$ | -0.01 | 0.01 | -0.03 | 0.01 |
| | *Group level effects* | | | | |
| | sd: Intercept | 0.88 | 0.41 | 0.34 | 1.64 |
| | sd: Baseline Slope | 0.06 | 0.05 | 0.00 | 0.15 |
| | sd: Level Change | 0.92 | 0.46 | 0.31 | 1.77 |
| | sd: Slope Change | 0.07 | 0.06 | 0.01 | 0.19 |
| | sd: Slope Change$^2$ | 0.01 | 0.01 | 0.00 | 0.03 |
| Unrelated | *Population level effects* | | | | |
| | Intercept | -3.46 | 0.66 | -4.59 | -2.45 |
| | Baseline Slope | -0.06 | 0.20 | -0.39 | 0.27 |
| | Level Change | 1.70 | 0.73 | 0.55 | 2.94 |
| | Slope Change | 0.51 | 0.24 | 0.12 | 0.91 |
| | Slope Change$^2$ | -0.04 | 0.02 | -0.07 | -0.02 |
| | *Group level effects* | | | | |
| | sd: Intercept | 0.79 | 0.51 | 0.11 | 1.71 |
| | sd: Baseline Slope | 0.10 | 0.08 | 0.01 | 0.24 |
| | sd: Level Change | 0.90 | 0.51 | 0.18 | 1.84 |
| | sd: Slope Change | 0.14 | 0.12 | 0.01 | 0.35 |
| | sd: Slope Change$^2$ | 0.02 | 0.01 | 0.01 | 0.05 |

**Supplementary materials S1: by-participant description of individualized BEARS training.**

**Participant 1** consistently benefited from cues to "Let the dust settle" (i.e., to allow for overactive lexical information in short-term memory to decay). This was because after multiple trials, he was noted to become perseverative, especially when he responded too quickly. Discussed "tool drawer" analogy (see Appendix 1). He was also observed to either build up too much tension prior to retrieving the target or move on too quickly without enough time for an attempt. Therefore, additional education was provided regarding balancing the right amount of time with the right amount of effort. He was cued to "breathe and slow down" as opposed to always saying "move on." Participant 1 also benefited from discussion about understanding and working *with* your aphasia. Discussed his apparent maladaptive habits of trying to force a word out through increased tension or trying to "sneak up on it" by trying to produce it as quickly as possible (which often resulted in errors), as opposed to a more adaptive approach of easing into word retrieval. Additional feedback regarding his level of initiation speed and apparent tension was provided as needed.

**Participant 2** had a tendency to make impulsive errors and benefited from cues to "breathe and slow down." Frequent tension feedback was provided when he became tense after missing multiple targets consecutively. Tension and frustration appeared to be his biggest battle. He benefited from discussion about understanding and working *with* your aphasia. Also discussed easing into retrieving a word rather than trying to "spit it out."

**Participant 3** benefited from education regarding how to modulate his naming speed. He was encouraged to slow down when he saw a picture that he knew gave him more difficulty. Discussed putting in the right amount of effort and being aware of his output (tool drawer analogy). He was encouraged to slow down if he found himself coming up with other treatment targets (i.e. visor for eggplant; both treated items). He was noted to become more aware of his output in the treatment game but not as much in conversation with the clinician. As a treatment session went on, he benefited from additional cues to breathe and relax, as tension would build up if he began missing treatment targets he could typically access. Of note, when increased difficulty was noted in a treatment session, it was frequently in conjunction with self-reports of poor sleep or stress in his personal life. Upon final follow up, participant's wife reported increased confidence. Participant reported that he was now ordering his own food at the deli counter without difficulty. Per wife, both of his children also reported noticeable improvements in his language.

**Participant 4** was good at providing his own life examples and analogies in relation to his experience of aphasia. He required a lot of cues and education to learn when to move on when a retrieval attempt was not going to be successful and when to give himself more time when he had a good chance of succeeding. Increased naming difficulty was noted during step 3 of each trial (after questions) compared to other participants. He was noted to be impulsive during his second attempts, frequently producing a word within the category but not the target item. He benefited from frequent cues to relax, as frustration was a barrier for his naming. Education was provided regarding the differences between intensive drill-based practice to improve his language system vs. typical total communication. He was encouraged to continue to use circumlocution strategies in his day-to-day life, but to only provide one word during the

treatment game. Upon exit, participant reported that he did not feel like he improved much in his speech/language but really enjoyed the treatment sessions.

**Participant 5** required multiple and frequent cues to provide a single word to name the item, as he provided constant commentary about his performance and was frequently tangential (e.g., sharing stories about his own life). Commentary and tangential stories were considerably reduced over the course of treatment. Frequent encouragement was provided to increase speed for semantic feature questions, as he often got distracted by commenting on the question or by responding quickly verbally but not selecting the button. Commentary was reduced with training, but speed did not seem to improve. He was encouraged to remain silent until he decided to attempt the word or move on. Noted improvement was made on this skill throughout treatment. He had a tendency to attempt to name a target word a number of times, which appeared to increase his chance of subsequent perseverative errors. He was instructed to name the picture only one time, and his ability to do so improved over the course of treatment. He appeared to benefit more from concrete examples and practice rather than from abstract analogies. He seemed to improve his ability to wait until he was ready before attempting to say the word. Tangential behavior (commentary and stories) were noted to increase with fatigue. Per wife, both his sister and brother-in-law commented on how much his language improved. Wife reported that his language was the best it has been since his stroke, and that he started to initiate more conversation in the car as opposed to sitting quietly.

**Participant 6** benefitted significantly from cues to "relax" (i.e. to self-monitor level of tension in her shoulders). She was encouraged to slow down, the difference between feeling rushed vs. not rushed as she attempted to retrieve a word. Overall, she was very motivated to earn feedback points and get through a many treatment trials. She was encouraged to self-monitor her output for sound errors. She was encouraged to use the "most efficient road" for word retrieval (natural access vs. compensatory strategies) during the treatment game. This was because she occasionally attempted to visualize and then read the spelling of a word in her head, which took longer and resulted in sound errors for irregularly spelled words. She was educated about modulating her speed as needed, and that she did not need to slow down for every picture, just when she encountered retrieval difficulty. She was also trained to only provide a single naming attempt once she was ready, as she often attempted to retrieve a word too quickly and then produced multiple self-corrected errors.

**Participant 7** was encouraged frequently to speed up his overall responses. He became extremely frustrated during more difficult tasks and would often take the maximum amount of time to respond. Used an education strategy related to his former law enforcement experience to illustrate the need to make a quick "gut decision" instead of seeking to "build a case" like a detective. Trained on ways to "clear the static" (i.e. taking a breath when the task became difficult, slowing down, counting to 10 to try to clear instances of perseveration). Encouraged him to be aware of output to reduce perseverations or unrelated responses. Both his wife and participant reported that his day-to-day language performance improved after treatment. Improved speed was evident in overall duration of exit testing.

**Participant 8** benefited from education regarding modulating his naming speed. He was encouraged to slow down when he saw a picture that he knew gave him more difficulty.

Discussed putting in the right amount of effort and self-monitoring output. He was encouraged to slow down if he found himself coming up with unrelated treatment targets (i.e. "visor" for "eggplant" which were both treated items). He implemented several strategies to relax and slow down without clinician prompting. Given comprehension impairments, he benefited most from training using direct strategies as opposed to abstract scenarios or analogies. He was very motivated by treatment game feedback. As treatment progressed and his performance accuracy increased, he was instructed to continue to balance speed and accuracy, performing at a faster rate that occasionally resulted in an error instead of slowing down significantly to achieve 100% performance accuracy.

**Participant 9** required significant encouragement throughout treatment. Following a few sessions of treatment, his increased awareness of errors and choosing to "move on" was observed to over-generalize to conversation, resulting in reduced output. Therefore, a distinction was drawn between the treatment game and general conversation, and he was encouraged to continue to focus on getting his point across despite sound errors (and to shift to alternative communication strategies, if necessary). Discussed putting in "the right amount" of effort and time before moving past a picture. Obvious tension was observed when he attempted a word before he was "ready." He was encouraged to "break the tension" by taking a deep breath and sitting back in his chair before trying again. With moderate cueing to complete these techniques, he had good success retrieving his target in instances of increased difficulty.

# Supplementary materials S2: Additional language testing results

## Additional participant testing: Philadelphia Naming Test

|  | Correct | Semantic Errors | Formal Errors | Mixed Errors | Unrelated Errors | Nonword Errors | Non-naming Response | S-Weight | P-Weight |
|---|---|---|---|---|---|---|---|---|---|
| Participant 1 | 72 | 23 | 6 | 4 | 6 | 35 | 29 | 0.0182 | 0.0619 |
| Participant 2 | 121 | 10 | 4 | 4 | 0 | 27 | 9 | 0.0294 | 0.0169 |
| Participant 3 | 151 | 6 | 2 | 5 | 0 | 6 | 5 | 0.0313 | 0.025 |
| Participant 4 | 129 | 18 | 5 | 12 | 1 | 5 | 5 | 0.02 | 0.0288 |
| Participant 5 | 63 | 8 | 9 | 8 | 3 | 7 | 77 | 0.0176 | 0.025 |
| Participant 6 | 117 | 0 | 13 | 1 | 1 | 30 | 13 | 0.0325 | 0.0151 |
| Participant 7 | 62 | 16 | 2 | 3 | 10 | 1 | 81 | 0.0151 | 0.0325 |
| Participant 8 | 103 | 6 | 17 | 3 | 3 | 23 | 20 | 0.0213 | 0.02 |
| Participant 9 | 120 | 5 | 14 | 3 | 1 | 22 | 10 | 0.0269 | 0.0182 |

**Additional participant testing: Nicholas and Brookshire, Camel and Cactus, PALPA 25 and PALPA 36**

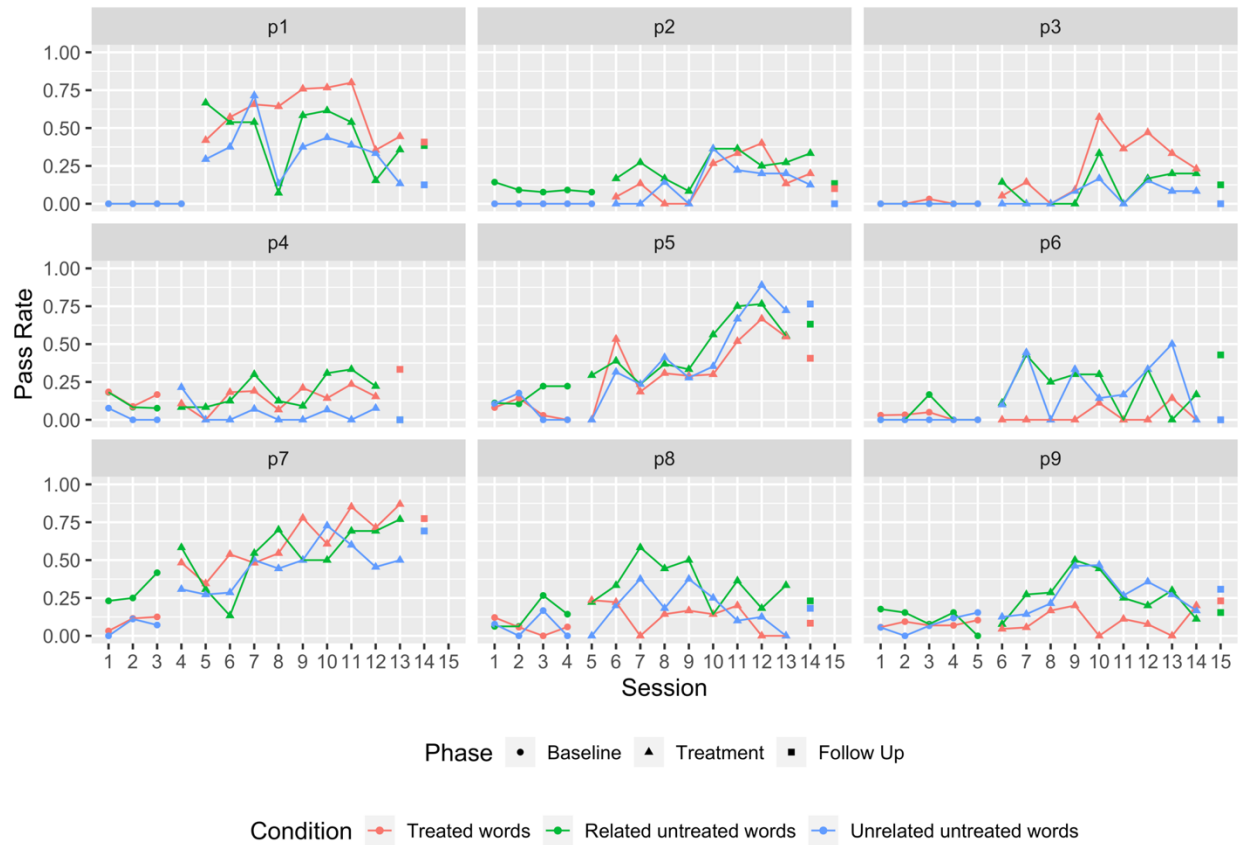| Assessments | Cactus and Camel | PALPA 25 | | PALPA 36 |
|---|---|---|---|---|
| | Percentage Correct | Words | Non-Words | |
| Participant 1 | 87.50% | 70% | 95% | 0 |
| Participant 2 | 75.00% | 100% | 83.33% | 7 |
| Participant 3 | 79.69% | 96.67% | 95% | 12 |
| Participant 4 | 71.89% | 78.33% | 96.67% | 0 |
| Participant 5 | 75.00% | 85% | 86.67% | 13 |
| Participant 6 | 92.19% | 100% | 96.67% | 21 |
| Participant 7 | 64.06% | 25% | 88.30% | 1 |
| Participant 8 | 84.40% | 93.33% | 98.33% | 18 |
| Participant 9 | 81.25% | 93.33% | 91.67% | 4 |

**Figure S3a.** Pass rate (proportion of "pass" responses relative to paraphasia and nonresponse error types) by session. P1, P5, and P7 were the three participants who responded least to the treatment in terms of naming probe accuracy.
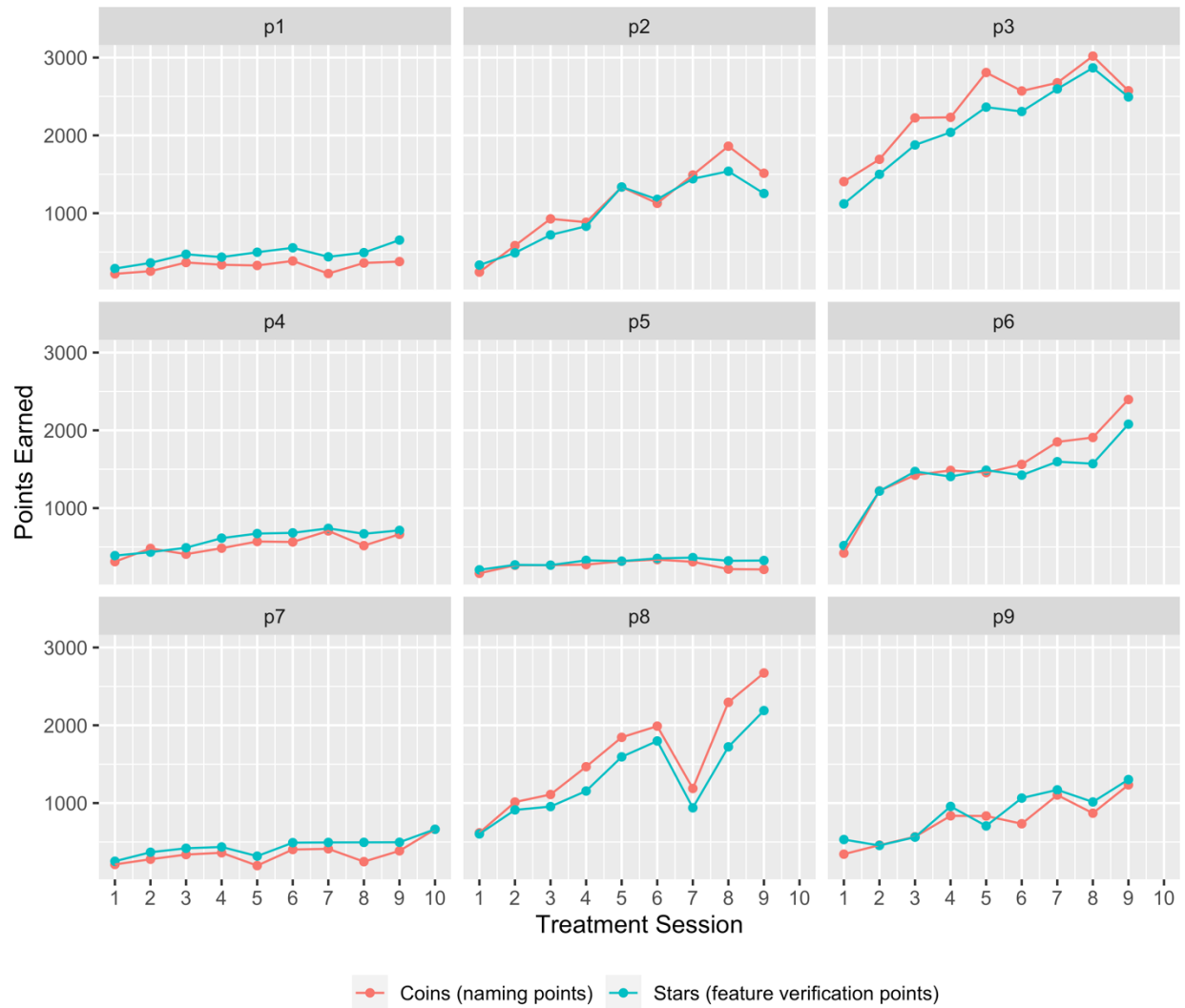
**Figure S3b.** Practice efficiency during treatment over time, as measured by cumulative feedback points by session.