

# A systematic review examining the efficacy of commercially available foreign language learning mobile apps

Jodi M. Tommerdahl, Chrystal Sapphire Dragonflame and Amanda A. Olsen

## QUERY SHEET

This page lists questions we have about your paper. The numbers displayed at left are hyperlinked to the location of the query in your paper.

The title and author names are listed on this sheet as they will be published, both on your paper and on the Table of Contents. Please review and ensure the information is correct and advise us if any changes need to be made. In addition, please review your paper as a whole for typographical and essential corrections.

Your PDF proof has been enabled so that you can comment on the proof directly using Adobe Acrobat. For further information on marking corrections using Acrobat, please visit <http://journalauthors.tandf.co.uk/production/acrobat.asp>; <https://authorservices.taylorandfrancis.com/how-to-correct-proofs-with-adobe/>

The CrossRef database ([www.crossref.org/](http://www.crossref.org/)) has been used to validate the references.

## AUTHOR QUERIES

- Q1** Please check and confirm whether the author affiliations and corresponding details have been set correctly.
- Q2** The keywords information has been imported from data supplied with the original manuscript. Please revise if incorrect.
- Q3** Please provide the short biography of the author in the “Notes on contributor” section.
- Q4** The disclosure statement has been inserted. Please correct if this is inaccurate.
- Q5** Please provide complete details for (ILTA, 2007, JLTA, 2002, EALTA, 2006, and Purpura, 2004) in the reference list or delete the citation from the text.
- Q6** There is no mention of (Clearinghouse 2012 and Purpura 2013) in the text. Please insert a citation in the text or delete the reference as appropriate.
- Q7** Please provide the volume number.

- Q8** Please provide the volume number.
- Q9** Please provide the year.
- Q10** Please provide the page range.
- Q11** Please note that the ORCID section has been created from information supplied with your manuscript submission/CATS. Please correct if this is inaccurate.

PROOF ONLY



# A systematic review examining the efficacy of commercially available foreign language learning mobile apps

Jodi M. Tommerdahl<sup>a,b</sup>, Chrystal Sapphire Dragonflame<sup>a,b</sup>  and Amanda A. Olsen<sup>a</sup> 

Q11

<sup>a</sup>Department of Curriculum and Instruction, University of Texas at Arlington, Arlington, TX, USA;  
<sup>b</sup>Southwest Center for Mind, Brain, and Education, University of Texas at Arlington, Arlington, TX, USA

## ABSTRACT

A systematic review examining the efficacy of commercially available foreign language-learning apps (FLL) was completed. A database search of ERIC, PsychINFO, and LearnTechLib produced 1,786 journal articles. After applying specific inclusion and exclusion criteria based on Burston's seminal study (2015) requiring a minimum number of 10 participants, quantitative learning outcome data and rigorous research design, eight studies remained. These studies were categorized in terms of the app studied, year of publication, language taught, age group of participants, setting, length of study, and device(s) used. Descriptive statistics demonstrate there is a dearth of studies examining app efficacy, that English was the most commonly taught language, and that vocabulary was the most commonly tested area. Although commercial apps were found to successfully support FLL, the included studies' methods varied in ways that made direct comparison difficult.

## KEYWORDS

Mobile apps; technology; foreign language-learning; mobile assisted language learning

Q2

Since the inception of Apple's app store and Android Market in 2008, mobile applications (apps) have grown exponentially in popularity, rapidly expanding to meet the enormous demand for novel varieties of apps (Pandey et al., 2019; Yu, 2019). For example, the Apple app store opened with 500 apps which by 2017 had expanded to 2.2 million. Although many apps, such as Angry Birds in 2009 (Rovio Entertainment Corporation, 2020), were created for their amusement value, others offered utility in areas such as social media, business, fitness, and education (Toto & Limone, 2019). Within the category of educational apps, the subcategory of foreign language learning (FLL) quickly emerged,

**CONTACT** Jodi M. Tommerdahl  [joditom@uta.edu](mailto:joditom@uta.edu)  Department of Curriculum and Instruction, University of Texas at Arlington, 701 S Nedderman Dr, Arlington, TX 76019, USA.

Q1

42 demonstrating the public's strong desire to learn foreign languages in  
43 an accessible and convenient way.

44 An early example of the enthusiasm and interest for FLL apps can  
45 be seen in the example of Babbel, which entered the market in 2007  
46 (Babbel, 2020). Although they do not release their overall subscriber  
47 numbers, they reported a gain of 50 million new downloads worldwide  
48 in 2018 alone (Chan, 2019).

49 Duolingo, arguably one of the most well-known mobile apps for  
50 language learning, also achieved rapid success. Within three weeks of  
51 launching in the Apple app store in 2012, the app reached one million  
52 downloads, with no advertising except by word of mouth. By 2013, four  
53 million individuals used the app, and in 2018 alone, the app was down-  
54 loaded another 50 million times (Chan, 2019). Duolingo offers over 30  
55 different languages to study with several more in either the development  
56 or beta stage (Duolingo, 2020). Although we are unsure of how many  
57 FLL apps exist in the world's app stores, a perusal of any major app  
58 download site allows us to estimate that at least several hundred are  
59 available, although the true number could be in the thousands.

60 Language learning apps vary widely from each other in qualities such  
61 as their selection of languages taught, and the native language used for  
62 instruction. For instance, an app developed for teaching Norwegian to  
63 Chinese speakers would not be appropriate for English speakers learning  
64 Norwegian. This means that some apps may have several different ver-  
65 sions for different audiences, even when teaching the same language.  
66 Apps also vary by the type of language skills they emphasize. Some  
67 apps may focus on understanding the pronunciation of the studied  
68 language, or concern themselves solely with vocabulary development,  
69 while others may attempt to blend several linguistic skills. Another  
70 variation between FLL apps is their revenue model. Revenue can be  
71 earned through advertising on the app, relying on in-app purchases,  
72 completing a one-time payment for access to the app, or through a  
73 recurring subscription fee.

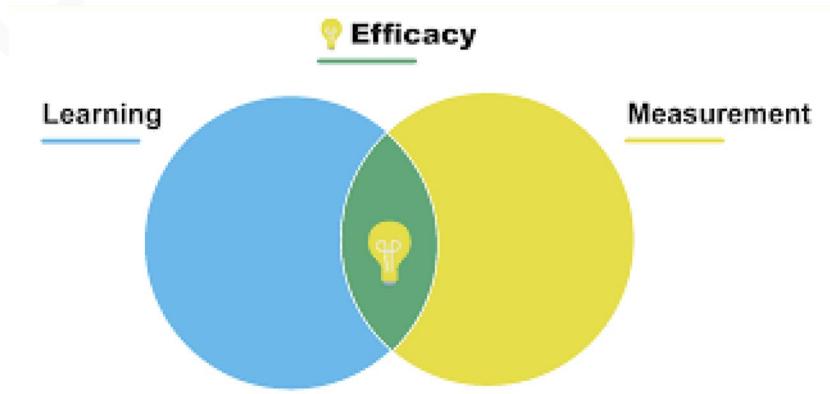
74 Given the large number of consumers who are using and often paying  
75 for FLL apps, the need for research into their efficacy is apparent, both  
76 for the consumer who is trying to decide between them and for app  
77 designers who want to ensure their apps are effective and attractive.  
78 The need for efficacy studies has been increasingly recognized in recent  
79 years in the field of education within core subject areas such as math-  
80 ematics (Doabler et al., 2019), science literacy (Goldman et al., 2019),  
81 and reading (Vaughn et al., 2019; Wanzek et al., 2011) with Toste et al.  
82 (2019, p. 46) referring to this need as "a central mission of educational  
83 research." This need is further evidenced by U.S. federal policy in the  
84 form of the Every Student Succeeds Act, a program which mandates

that federal money must be used for programs that show evidence of effectiveness (Dynarski, 2015) and the creation of the What Works Clearinghouse, a website that provides research-based efficacy guidelines and information in several areas of education (<https://ies.ed.gov/ncee/wwc/>). When considering efficacy, it is important to define both efficacy and an efficacy perspective. The current paper draws heavily from the Efficacy Framework put forward by American College Testing (ACT) (Mattern, 2019) which defines efficacy as “the degree to which evidence, rationales, and theory support the claim that a learning tool improves intended learner outcomes under ideal conditions” (p. 2). More specifically, the framework targets ‘results/impact,’ defined as “the degree to which targeted outcomes occur as a result of using the learning tool.” This is important as designers of instructional tools need to base their work on specific knowledge components to optimize learning (Mattern, 2019).

The Efficacy Framework seeks to align the validity of assessments and learning product efficacy, allowing them to achieve two stated goals:

1. Develop learning solutions that are most likely to impact the intended outcome.
2. Increase our ability to detect whether a learning solution is achieving its intended outcome through thoughtful study design, methodology, and data collection.

The ACT Framework proposes that the quantification of efficacy arises from the intersection of learning and measurement, as shown in Figure 1. This requires measuring the impact of the learning tool on the intended outcome, which is in turn aligned with the Knowledge-Learning-Instruction Framework (Koedinger et al., 2012).



**Figure 1.** ACT’s efficacy framework: intersection of learning and measurement.

128 To judge the efficacy of instructional tools, apps in this case, it is  
129 important that rigorous efficacy studies are completed. Then systematic  
130 reviews can help compare and contrast the varying studies to accurately  
131 evaluate the research. Although several reviews of mobile assisted lan-  
132 guage learning (MALL) exist, none to our knowledge focus exclusively  
133 on commercially available mobile phone and tablet apps designed for  
134 FLL. This study's objective was to focus on this specific area, firstly to  
135 determine the amount of rigorous research into investigating their effi-  
136 cacy, and secondly, to examine the results of those publications in an  
137 effort to help individuals such as language learners, foreign language  
138 instructors, and app designers to become more aware of which apps  
139 have been shown to be effective.  
140

### 141 **Literature review**

143 Although no efficacy studies of FLL mobile apps appeared in our search  
144 of existing literature, FLL mobile apps have been studied in other con-  
145 texts apart from their efficacy such as within the larger categories and  
146 MALL. The following literature review reports on these related areas  
147 not exclusive to FLL mobile app efficacy, but which nonetheless report  
148 on FLL apps.  
149

### 150 ***Review of FLL: Games without efficacy measures***

152 In a FLL review, Dehghanzadeh et al. (2019) studied the effect of games,  
153 many of which were apps, on learning English as an additional language  
154 in a digital environment. Twenty-two studies were reviewed, and all  
155 reported positive experiences for students using ESL games. However,  
156 the efficacy of language learning was not measured. Instead, the study  
157 examined engagement and motivation. The contents of the 17 of 22  
158 games studied were mostly focused on vocabulary with five focusing  
159 on grammar, four on pronunciation, five on speaking, four on listening,  
160 and three on writing. The overall finding was that participants found  
161 the games to be “enjoyable, engaging, motivating, and fun” (Dehghanzadeh  
162 et al., 2019, p. 1).  
163

### 164 ***Review specific to apps for FLL but not measuring efficacy***

166 In a similar paper, Heil et al. (2016) reviewed the 50 most popular FLL  
167 commercial apps with regard to three questions: 1) What are the primary  
168 pedagogical focuses of popular language learning apps?; 2) Do apps  
169 adapt to individual needs, language proficiency levels, and styles of  
170 learning?; and 3) How is corrective feedback employed in these apps?

171 These questions stemmed from theoretical models of language which  
172 categorized different aspects of language learning such as knowledge of  
173 grammar, pragmatics, discourse, and sociolinguistics (Bachman & Palmer,  
174 1996; Purpura, 2004). Further, Heil et al. (2016) argued that several  
175 different areas of language must be integrated for effective language  
176 learning, rather than promoting the sheer memorization of words. Their  
177 review found that 42 out of 50 commercial FLL apps taught vocabulary  
178 in isolation, whereas only 12 offered grammar instruction, demonstrating  
179 that a majority of the most popular commercial FLL apps were not  
180 appropriately equipped for teaching conversational language. This sup-  
181 ports Dehghanzadeh et al. (2019) results which discovered an emphasis  
182 on vocabulary in language learning games.

### 184 ***Reviews measuring efficacy of MALL for FLL but not limited to apps***

186 Although previous reviews have analyzed different aspects of technology  
187 and language learning, few have examined the actual efficacy of tech-  
188 nological tools to enhance learning (Bolgün & McCaw, 2019). This was  
189 surprising, given the emphasis that researchers and government bodies  
190 have placed on understanding the efficacy of educational interventions.  
191 Furthermore, consumers who are using, and often paying for FLL apps,  
192 with the intent of learning a new language quickly and efficiently, should  
193 have research-based guidance available when making their choice.  
194 Although some data regarding the efficacy of FLL apps exists, it is  
195 typically embedded in studies that have analyzed the larger category of  
196 mobile assisted language learning (MALL). A discussion of these studies  
197 follows, including descriptions of their methodologies since they are  
198 meta-analyses, as is the current paper. Note, each of the following MALL  
199 studies included mobile apps, but their data was not analyzed separately  
200 from the other types of MALL, thereby not allowing readers to measure  
201 the efficacy of mobile apps specifically.

202 For example, Sung et al. (2015) completed a meta-analysis analyzing  
203 the efficacy of mobile devices for teaching FLL, extending beyond apps  
204 to include items such as text messaging, social media, global positioning  
205 systems, and video capture. They also examined commonalities found  
206 within MALL articles published between 1993 and 2013, a timeframe  
207 notably beginning before mobile apps, by examining the type of par-  
208 ticipants, hardware and software, teaching and learning methods, settings,  
209 language skills, target languages, and intervention duration. From an  
210 initial search returning 721 results including journal articles, conference  
211 papers, and doctoral dissertations, two criteria were applied: 1) the  
212 research question was required to be about mobile device use in foreign  
213 language learning, and 2) only experimental and quasi-experimental

214 studies that allowed for the calculation of effect sizes were included.  
215 Screening using these criteria reduced the final number of articles to  
216 44. Their results reported that the most common age group of partic-  
217 ipants was elementary-school students, the most common timeframe for  
218 experiments using learning-oriented software was one to six months,  
219 the most frequently studied learning skill was vocabulary, and the most  
220 common language being learned was English.

221 When analyzed statistically, 70% of participants using mobile devices  
222 for FLL outperformed their control group, providing evidence of FLL  
223 mobile device effectiveness. Further findings revealed that adults and  
224 school children had similar results from using MALL, handhelds had a  
225 larger effect size than laptops, use of mobile devices in multiple learning  
226 settings was more effective than the more restricted settings of the  
227 outdoors or the classroom, interventions of 1-6 months were more effec-  
228 tive than shorter or longer time periods, and finally, using mobile devices  
229 for vocabulary or mixed skills produced higher achievement than those  
230 which focused on individual skills such as listening or reading. Despite  
231 the value of Sung et al. (2015) findings, the research does not provide  
232 data specific to mobile apps although it includes technologies such as  
233 social media and video capture.

234 Additionally, Burston (2015) completed a rigorous systematic review  
235 of the learning outcomes of MALL implementation projects. This  
236 included 291 articles published between 1994 and 2012, again mainly  
237 predating apps, but examining a range of technologies including video  
238 playback, flashcards, and speech recordings. Inclusion criteria applied  
239 to these 291 articles included having a minimum of 10 participants  
240 using the mobile device for a minimum of one-month, quantitative  
241 learning outcome data, tracing the amount of usage of the mobile device,  
242 adequate control group descriptions, absence of uncontrolled and con-  
243 founding variables, and adequate statistical analyses. After excluding  
244 papers that did not meet this threshold, only 19 of the 291 studies could  
245 meaningfully evaluate the efficacy of MALL usage.

246 Out of the 19 studies, children and adult participants were equally  
247 represented with vocabulary being the most studied area. The most  
248 common treatment duration was from four to six weeks with the most  
249 common number of participants being in the 25- to 49-person range.  
250 An attempt to calculate effect sizes to compare across papers resulted  
251 in nine more papers being excluded due to their failure to report the  
252 language level of the participants. The final 10 papers were split between  
253 five different areas of language, vocabulary, reading, speaking, listening,  
254 and writing, making a comparison of effect sizes impossible. Despite  
255 the large number of MALL papers published, this review draws attention  
256 to how few reliable studies were available regarding MALL efficacy. It

257 is particularly interesting to notice that these two reviews restricted  
258 inclusion criteria to quantitative, experimental and quasi-experimental  
259 studies, therefore, the initial number of studies found decreased greatly,  
260 showing that most research was limited in its ability to provide quality  
261 statistical results of efficacy. Sung et al. (2015) research included only  
262 6.1% of the initial articles found while Burston's (2015) included only  
263 3.4%, likely due to the stricter criteria including number of participants,  
264 minimum amount of time of study length, and quality statistical analyses  
265 among others.

266 The limited amount of research on the efficacy of MALL, long pre-  
267 dating 2008, emphasizes the question of how much research has been  
268 completed on the efficacy of FLL apps that are available to consumers  
269 in app stores. While FLL apps have been included in larger studies of  
270 efficacy, they have not been examined independently as a group. Although  
271 commercial FLL apps generally claim to use effective pedagogical concepts  
272 in their design (Toto & Limone, 2019), it is unclear whether research  
273 evidence to support their efficacy exists, and if so, for which apps and  
274 for what language areas. Without this knowledge, it is impossible for  
275 consumers to make informed decisions on whether using an app for FLL  
276 is beneficial, and if so, which are the most effective. The intent of the  
277 present systematic review was to answer the following questions:

- 278 1. How many scientifically reputable studies exist on FLL app efficacy  
279 showing learning outcomes?
- 280 2. What general trends exist in the research (number of studies per  
281 year and age groups studied)?
- 282 3. What languages do they teach and to what linguistic audiences?
- 283 4. What specific linguistic skills have been measured in these  
284 studies?
- 285 5. What are the efficacy outcomes of the studies?
- 286
- 287

## 288 **Methods**

289

290 This study was part of a larger study examining the efficacy of all FLL  
291 apps, whether commercially available or not (Olsen et al., under review).  
292 Although a close examination of FLL apps that are not commercially  
293 available may be of interest for professional app developers, it was felt  
294 that the public's need for information would be better served by focusing  
295 on obtainable apps. Similarly, language teachers looking for mobile apps  
296 to assist their students may be interested in the efficacy of commercial  
297 apps compared to those that are unpublished and unavailable.

298 To answer the five research questions, a systematic review following  
299 the Preferred Reporting Items for Systematic Reviews and Meta-Analyses

300 for Systematic Reviews Protocols (PRISMA-P) (Moher et al., 2015;  
301 Shamseer et al., 2015) was completed in March of 2020. Moher et al.  
302 (2015) explained that systematic reviews should “build on a protocol  
303 that describes the rationale, hypothesis, and planned methods of the  
304 review” (p.1). The PRISMA-P is a list consisting of 17 items ensuring  
305 that reviews are prepared and reported in a rigorous and robust manner  
306 and is often used in healthcare. A systematic review was chosen for  
307 this study since studies may have different outcomes, used different  
308 participants, and/or different mobile apps, meaning a meta-analysis  
309 would not be appropriate (Sriganesh et al., 2016).

310 After carrying out an extensive search for existing systematic reviews  
311 on the efficacy of FLL apps and finding none to exist, a list of criteria  
312 was developed, based on Burston’s (2015) article. These criteria are  
313 listed below.

#### 314 *Criteria for Inclusion*

- 315
- 316 1. Published in peer-reviewed journals from 2008 onwards
- 317 2. Written in English
- 318 3. Stated an intention to examine the efficacy of an app in relation  
319 to FLL outcomes
- 320 4. Implemented a research methodology (qualitative, quantitative or  
321 mixed methods)
- 322 5. Used mobile apps that were designed to teach a foreign  
323 language
- 324 6. Used mobile apps that were commercially available in an app store
- 325

#### 326 *Criteria for Exclusion*

- 327
- 328 1. Fewer than 10 participants in the study
- 329 2. Research design shortcomings such as the existence of uncontrolled  
330 variables, lack of a control group or inadequate statistical  
331 analyses
- 332

333 Although most criteria were identical to Burston’s, two exclusion  
334 criteria were dropped and one inclusion criterion was added. Specifically,  
335 Burston required the amount of time interacting with an app to be a  
336 minimum of one month; however, it was decided that an app showing  
337 efficacy in language learning in a shorter period of time would be of  
338 interest to the public. Burston’s limitation of articles to those that  
339 included an in-app tracking of time was also dropped, as the authors  
340 felt that the amount of time an app was open on a participant’s mobile  
341 device did not necessarily reflect focus and attention given to the app.  
342 The decision to omit two of Burston’s criteria allowed us the potential  
to include more studies. Regarding inclusion criteria, this study added

343 the requirement that apps studied be commercially available in an app  
344 store, thereby allowing the public and FLL educators to make decisions  
345 regarding apps that are actually available to them instead of having to  
346 sort through numerous studies where apps were designed specifically  
347 by researchers for a study and were not available to the public.

348 With the assistance of a university librarian, strings of search terms  
349 were developed and entered into the databases. The specific search  
350 strings used are provided in [Table 1](#).

351 Databases searched included the Educational Resources Information  
352 Center (ERIC), PsychINFO, and LearnTechLib. This search produced  
353 1,786 studies which were uploaded into Covidence, a software specifically  
354 designed for the development of systematic reviews. Using this software  
355 and the rules of PRISMA-P, the inclusion and exclusion criteria were  
356 applied to each article by the first (A.O.) and third (J.T.) authors. Reviews  
357 were completed in two stages, the abstract review, followed by a full-text  
358 review when appropriate. At each stage the author accepted or rejected  
359 each paper blindly, without the input of the other author until all articles  
360 had been reviewed and categorized. Only at that point were conflicts  
361 addressed. These conflicts were resolved through the close examination  
362 of each article and discussion between the authors. Each remaining  
363 article was categorized according to several variables of interest by the  
364 final author (J.T.). The second author (C.D.) also re-categorized 10% of  
365 the articles that were randomly selected in order to ensure reliability.  
366 Agreement across the variables of interest was over 95%.

## 367 **Results**

368  
369  
370 The first and most striking result was that only eight studies met the  
371 inclusion criteria, even with the easing of some of Burston's restrictions  
372 (see [Figure 2](#)).

373 This stands in stark contrast to what was initially expected given that  
374 a search of "Duolingo" on academic databases shows numerous results.  
375 Although it appears that a great deal is being published about FLL apps,  
376 very few are rigorous efficacy studies interested in learning outcomes.  
377 [Figures 3–13](#) display general trends of this limited group of studies.

378 None of the papers in this review were published until 2016, eight  
379 years after the introduction of FLL apps to app stores. Research of this  
380 type does not appear to be increasing since that time, with an average  
381 of approximately two papers a year since 2016. No developed body of  
382 work in peer reviewed journals exists with regard to any individual FLL  
383 app. In fact, the only app to have been effectively studied for efficacy  
384 in more than one paper was Duolingo which was represented here only  
385 twice (25%).

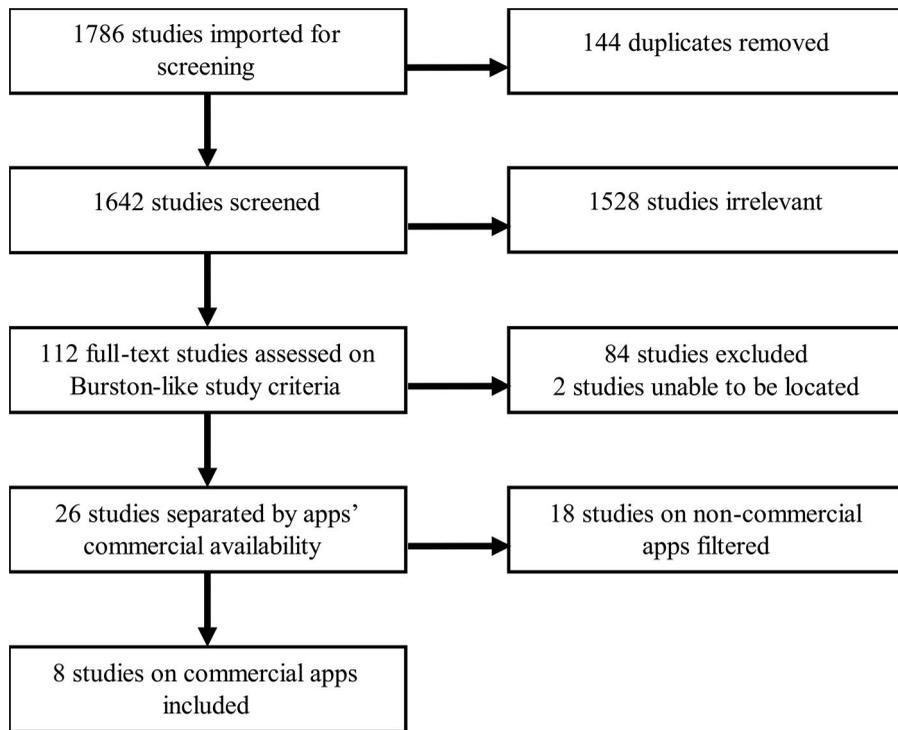
**Table 1.** The protocol executed in each database.

| Search Order | Database  | Protocol   |
|--------------|-----------|--|
| 1            | ERIC      | TX (mobile applications or apps or smartphone or iphone or iOS or android or lingualift or duolingo or hellotalk or mindsnacks or busuu or babbel or triplingo or mosalingua or cell phones or ipods or mobile devices or mobile phones or memrise or hinative or "(How to) Pronounce" or mondiy or lirica or drops or pimsleur or beelinguapp or clozemaster or fluento or acellastudy essential apps or tandem or lyrics training or "mobile-assisted language learning" or MALL) or [DE]Handheld devices  |
| 2            | ERIC      | ( (DE "Second Language Learning" OR DE "Bilingual Education" OR DE "Bilingual Instructional Materials" OR DE "Bilingualism" OR DE "College Second Language Programs" OR DE "Communicative Competence (Languages)" OR DE "English (Second Language)" OR DE "Language Enrichment" OR DE "Second Language Instruction" OR DE "Second Language Programs") OR TX "second language acquisition" or TX "mobile-assisted language learning"  |
| 3            | ERIC      | (DE "Academic Achievement" OR DE "Educational Attainment" OR DE "Student Promotion" OR DE "Academic Ability" OR DE "Academic Aptitude" OR DE "Academic Aspiration" OR DE "Academic Failure" OR DE "Academic Probation" OR DE "Achievement Rating" OR DE "College Readiness" OR DE "Educational Indicators" OR DE "Excellence in Education" OR DE "Grades (Scholastic)" OR DE "Grading" OR DE "Growth Models" OR DE "Instructional Effectiveness" OR DE "Knowledge Level" OR DE "Performance" OR DE "Reading Achievement" OR DE "School Effectiveness" OR DE "Student Evaluation" OR DE "Vocabulary Development" OR DE "Verbal Development" OR DE "Language Acquisition" OR DE "tests" OR DE "Scores" OR DE "Success" DE "Performance Based Assessment" OR DE "Language Aptitude" OR DE "Language Fluency" OR DE "Language Proficiency" OR DE "Language Skills" OR DE "Language Tests") OR TX (student success or GPA or grade point average or academic performance or learning performance or test scores OR vocabulary learning or vocabulary acquisition) |
| 4            | ERIC      | #1 AND #2 AND #3   |
| 5            | ERIC      | #5 Limits: English, Scholarly (Peer Reviewed) Journals, Academic Journals, Publication Date 2008-  |
| 1            | PsychINFO | #1 TX (apps or iphone or iOS or android or lingualift or duolingo or hellotalk or mindsnacks or busuu or babbel or triplingo or mosalingua or cell phones or ipods or mobile devices or mobile phones or memrise or hinative or "(How to) Pronounce" or mondiy or lirica or drops or pimsleur or beelinguapp or clozemaster or fluento or acellastudy essential apps or tandem or lyrics training or "mobile-assisted language learning" or MALL) or DE "Smartphones" OR DE "Mobile Phones" OR DE "Digital Gaming" OR DE "Mobile Applications" OR DE "Smartphone Use"  |
| 2            | PsychINFO | #2 DE "Foreign Language Learning" OR DE "Bilingual Education" OR DE "Foreign Languages" OR DE "English as Second Language" OR DE "Language Development" OR DE "Language Proficiency" OR DE "Lexical Access" OR DE "Linguistics" OR DE "Foreign Language Education" OR DE "Vocabulary" OR TX (second language learning or second language acquisition or "mobile-assisted language learning")   |
| 3            | PsychINFO | DE "Academic Achievement" OR DE "Academic Overachievement" OR DE "Academic Underachievement" OR DE "College Academic Achievement" OR DE "Academic Achievement Motivation" OR DE "Academic Achievement Prediction" OR DE "Academic Aptitude" OR DE "Academic Failure" OR DE "Verbal Comprehension" OR DE "Listening Comprehension" OR DE "Reading Comprehension" OR DE "Sentence Comprehension" OR TX (student success or GPA or grade point average or academic performance or learning performance or test scores OR vocabulary learning or vocabulary acquisition)   |
| 4            | PsychINFO | #1 AND #2 AND #3   |
| 5            | PsychINFO | Limits: English, Scholarly (Peer Reviewed) Journals, Academic Journals, Publication Date 2008-   |

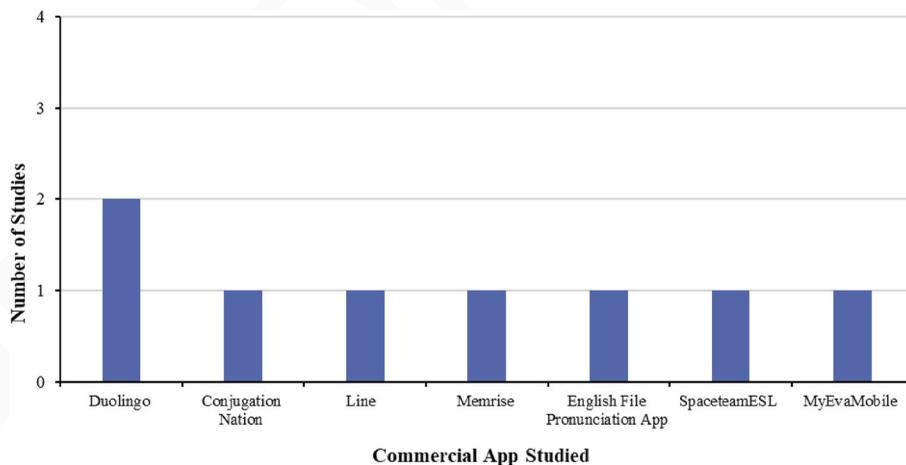
(Continued)

**Table 1.** Continued.

| Search Order | Database                                | Protocol   |
|--------------|---|--|
| 1            | LearnTechLib (full texts and abstracts) | (mobile applications OR apps OR smartphone OR iphone OR iOS OR android OR lingualift OR duolingo OR helloworld OR mindsnacks OR busuu OR babbel OR triplingo OR mosalingua OR cell phones OR ipods OR mobile devices OR mobile phones OR memrise OR hlnative OR "(How to) Pronounce" OR mondiy OR lirica OR drops OR pimsleur OR beelinguapp OR clozemaster OR fluento OR accelastudy essential apps OR tandem OR lyrics training OR "mobile-assisted language learning" OR MALL OR handheld devices) AND (foreign language learning OR bilingual education OR foreign language OR english as second language OR language development OR language proficiency OR lexical access OR linguistics OR foreign language education OR vocabulary) AND (academic achievement OR academic overachievement OR academic ability OR academic aptitude OR academic aspiration OR academic failure OR academic probation OR educational indicators OR academic excellence OR grades OR grading OR instructional effectiveness OR knowledge level OR performance OR reading achievement OR student evaluation OR vocabulary development OR verbal development OR language acquisition OR tests OR scores OR success OR language aptitude OR language fluency OR language proficiency OR language skills OR language tests OR student success OR GPA OR grade point average OR academic performance OR learning performance OR vocabulary learning OR vocabulary acquisition)   |
| 2            | LearnTechLib (journal articles)         | (MAINSUBJECT("Cell Phones") OR "smartphones" OR "tablet computers" OR "portable computers" OR "handheld computers" OR "apps" OR "iphone" OR "iOS" OR "android" OR "lingualift" OR "duolingo" OR "helloworld" OR "mindsnacks" OR "busuu" OR "babbel" OR "triplingo" OR "mosalingua" OR "cell phones" OR "ipods" OR "mobile devices" OR "mobile phones" OR "memrise" OR "hlnative" OR "(How to) Pronounce" OR "mondiy" OR "lirica" OR "drops" OR "pimsleur" OR "beelinguapp" OR "clozemaster" OR "fluento" OR "accelastudy essential apps" OR "tandem" OR "lyrics training" OR "mobile-assisted language learning" OR "MALL") AND (MAINSUBJECT("Second Language Learning") OR MAINSUBJECT("Bilingual Education") OR MAINSUBJECT("Second Languages") OR MAINSUBJECT("Linguistic Interference") OR MAINSUBJECT("Direct Method of Language Teaching") OR MAINSUBJECT("Audiolingual Language Teaching") OR "second language learning" OR "second language acquisition" OR "mobile-assisted language learning" OR "MALL" OR "language teaching methods" OR "language acquisition" OR "foreign language learning") AND (MAINSUBJECT("Learning") OR MAINSUBJECT("Achievement Tests") OR MAINSUBJECT("Academic Achievement") OR MAINSUBJECT("Tests") OR MAINSUBJECT("Reading Comprehension") OR MAINSUBJECT("Reading Ability") OR MAINSUBJECT("Reading Tests") OR MAINSUBJECT("Reading Acquisition") OR MAINSUBJECT("Reading Achievement") OR MAINSUBJECT("Verbal Learning") OR MAINSUBJECT("Listening Comprehension") OR MAINSUBJECT("Language Tests") OR MAINSUBJECT("Pronunciation Accuracy") OR MAINSUBJECT("Language Proficiency") OR MAINSUBJECT("English Proficiency") OR "student success" OR GPA OR "grade point average" OR "academic performance" OR "learning performance" OR "test scores" OR "vocabulary learning" OR "vocabulary acquisition" OR "language proficiency" OR "pronunciation accuracy" OR "reading achievement") AND stype:exact("Scholarly Journals") AND at:exact("Article") AND la:exact("English") AND PEER(yes) |



412 **Figure 2.** PRISMA diagram for foreign language mobile apps published between 2008 and  
413 2020.



428 **Figure 3.** Number of studies per commercial app.

430  
431  
432  
433  
434

The number of target languages for learning was limited to three with only one French (12.5%), two Spanish (25%), and five English (62.5%) papers represented. The language that apps used to communicate with the users were French, Spanish, English, Persian, and Chinese.

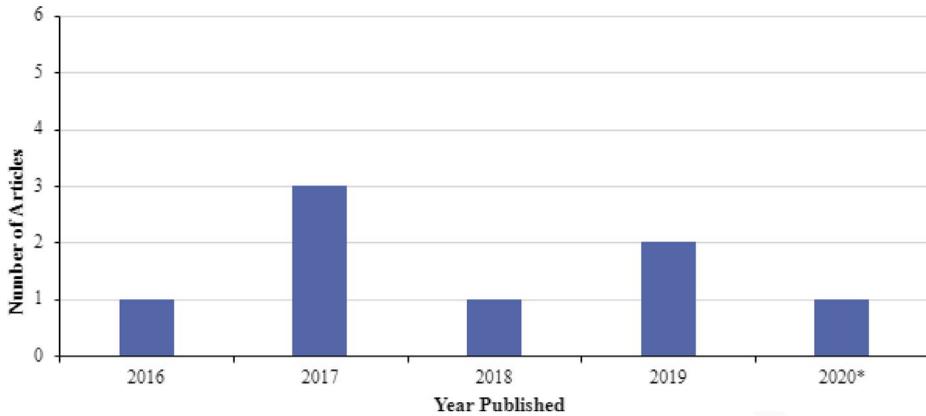


Figure 4. Number of studies published per year.

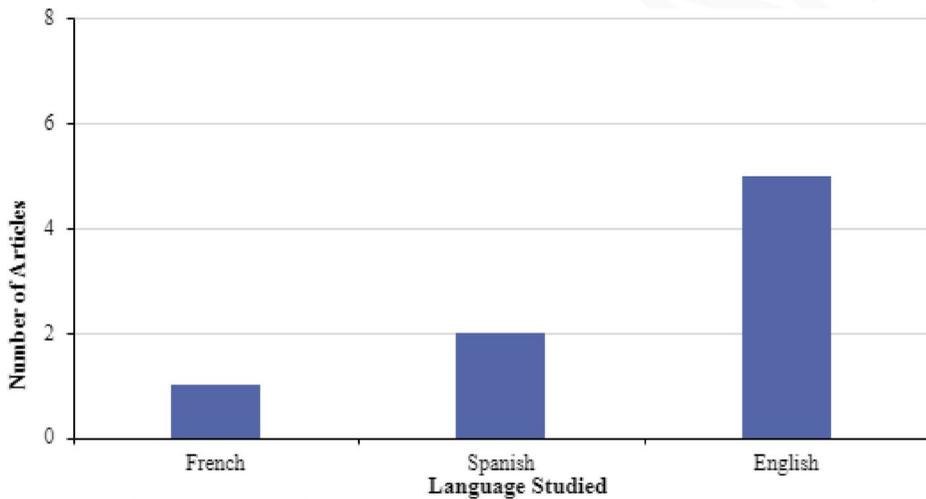


Figure 5. Target language studied.

Countries where studies were completed covered several nations, with seven countries producing the eight studies, two of which occurred in the U.S. (25%). The specific linguistic skills measured varied between studies, with three (37.5%) focusing on vocabulary, one (12.5%) on pronunciation, one (12.5%) on more than two skills, and one each for grammar (12.5%), oral fluency (12.5%), and 'listening and reading' (12.5%). Participants in seven (87.5%) of those studies were college students and app usage was completed solely in class for five (62.5%) of the eight studies. The length of studies varied from two weeks to four consecutive semesters. If the study had followed the four-week minimum used by Burston (2015), the number of articles included would have dropped to six (75%). It appears that mobile phones were the main

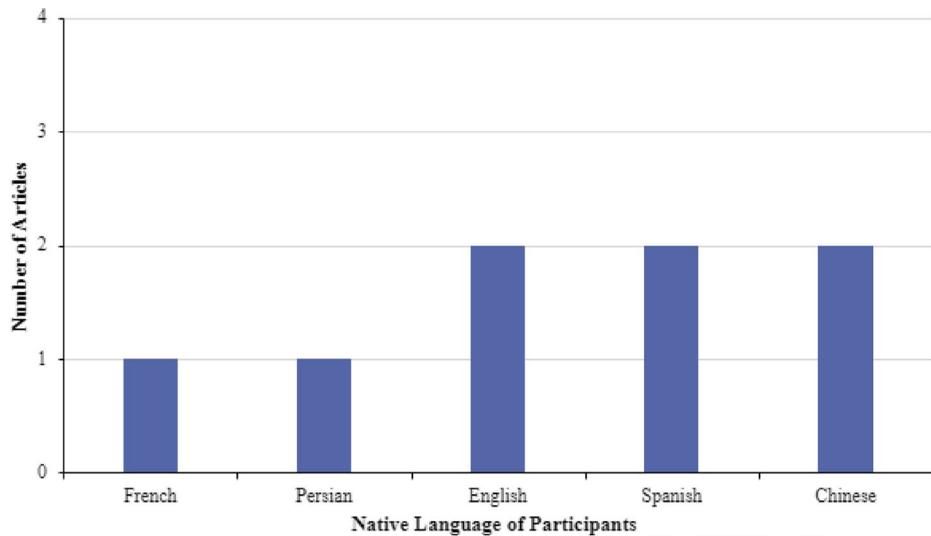


Figure 6. Native language of participants.

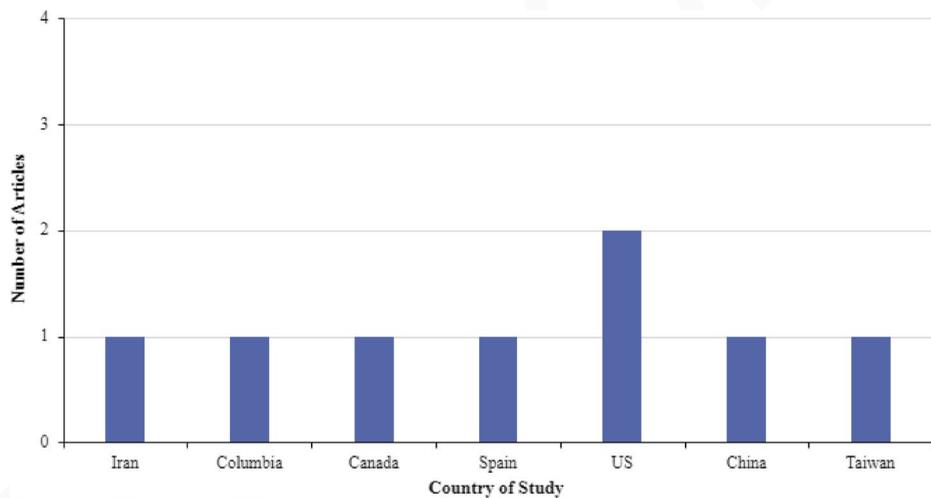
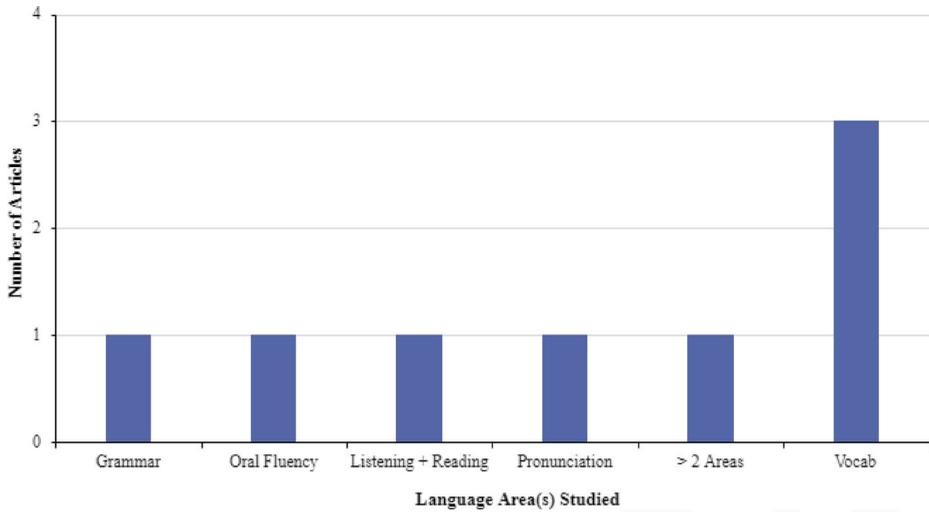


Figure 7. Country where research was carried out.

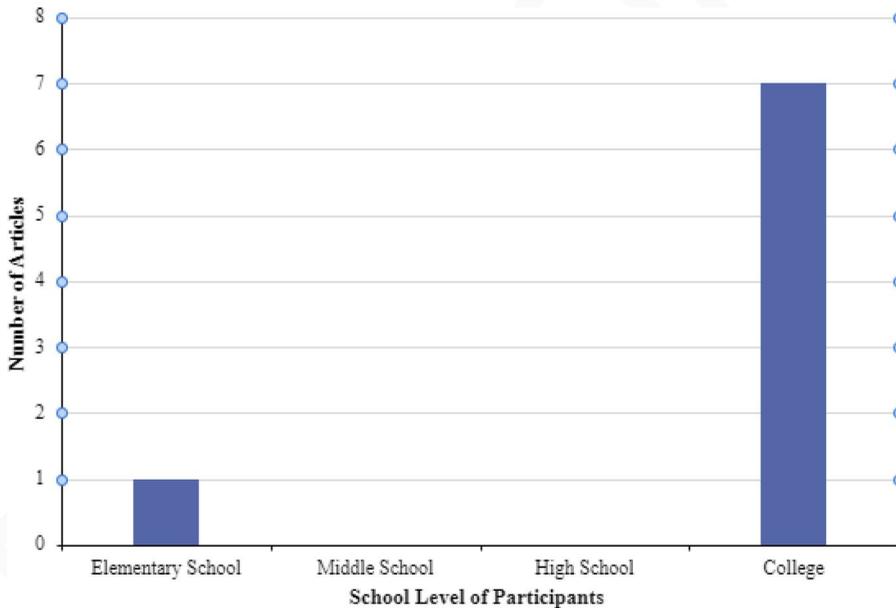
device used compared to tablets, but this is not certain as some studies did not specify the tool used.

With only eight studies meeting the systematic review's criteria, it is possible to present them individually to help better understand the differences and similarities between them. Their results are in [Table 2](#).

García Botero et al. (2019) Duolingo study divided 52 Columbian college students into two experimental groups and one control group. The control group did not use the app and received traditional classroom instruction only. Both experimental groups were provided with an introduction to the app in addition to traditional classroom instruction, but



535 **Figure 8.** Language area(s) studied.



554 **Figure 9.** School level of participants.

555  
556  
557  
558  
559  
560  
561  
562  
563

they varied in that one experimental group, the 'self-regulation and scaffolding group' received instruction about self-regulation strategies. This group was supported for self-regulation through weekly subgoals and scaffolding by the instructor and/or the service desk. After ten weeks, no significant differences were found between the control and experimental groups on listening, reading, or writing, ( $p > .05$ ). Students in the experimental group with self-regulation and scaffolding performed better in writing than the regular experimental group ( $p < .05$ ).

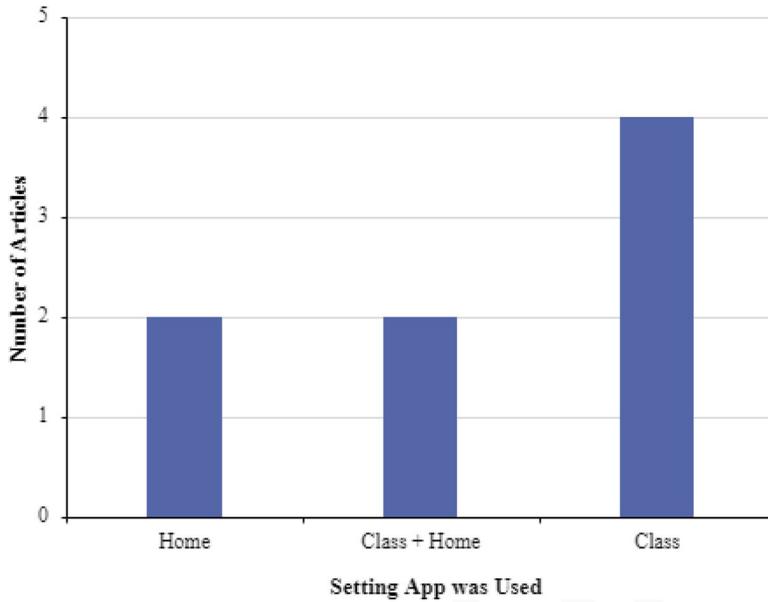


Figure 10. Setting of app use.

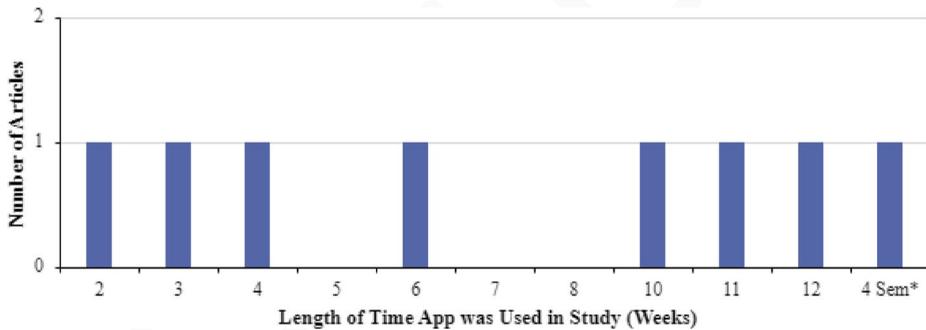
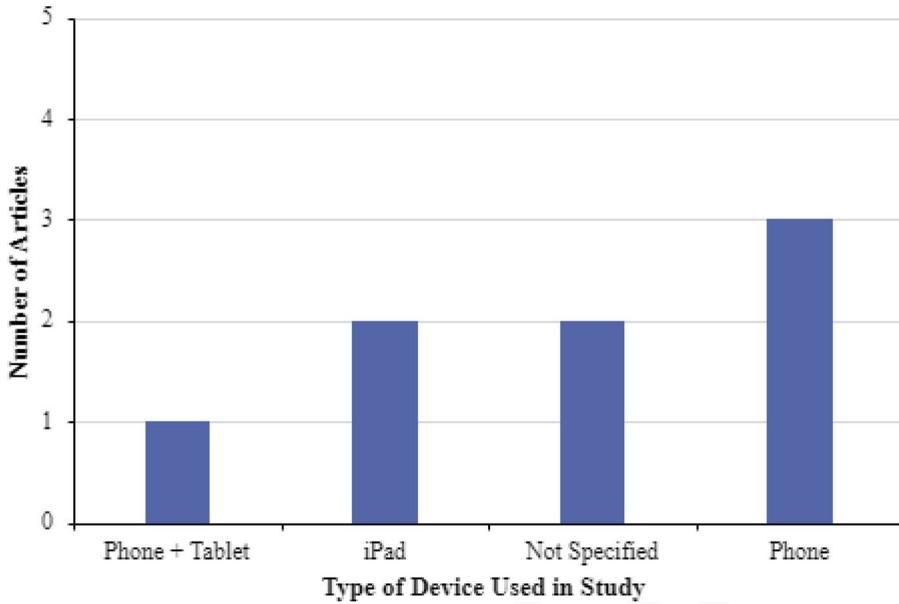


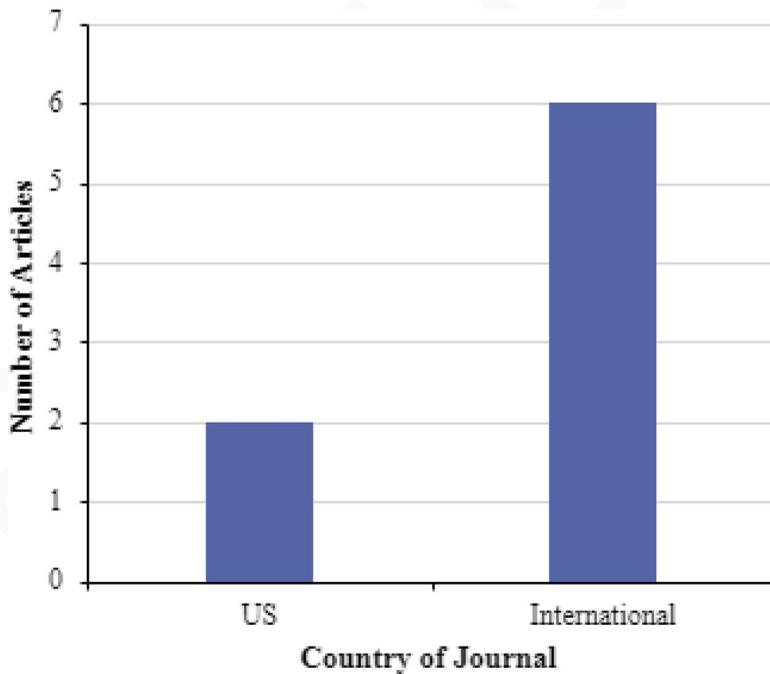
Figure 11. Studies' length of time.

The other Duolingo study (Rachels & Rockinson-Szapkiw, 2018) divided 164 US-based third and fourth grade students into an experimental and control group. For a period of 12 weeks, the experimental group only received Duolingo and the control group received traditional face-to-face teaching methods. Note the curriculum for the control group was adapted to mirror the content of Duolingo. Language knowledge was tested with a researcher-created tool that had vocabulary as its focus with the addition of some grammar-based test items. The results showed there were no statistically significant differences between the experimental and control group on the test ( $p > .05$ ).

Castañeda and Cho (2016) examined the app Conjugation Nation with 80 U.S. university students enrolled in Spanish courses designed to



624 **Figure 12.** Type of device used in study.



644 **Figure 13.** Country of publishing journal.

645  
646  
647 specifically measure improvement in verb conjugation ability. The app  
648 was used by students in pairs or trios during their classes where they  
649 also received traditional instruction for four consecutive semesters.

**Table 2.** Included studies by variables of interest.

| Author (Year)                      | App Name                                 | Skill(s) Measured                                 | Country of Study         | Language Studied | Outcome  |
|------------------------------------|--|---|--------------------------|------------------|--|
| García Botero et al. (2019)        | Duolingo                                 | Listening, reading, writing                       | Columbia                 | French           | No significant differences were found between the control group and experimental groups in any area. |
| Rachels & Rockinson-Szapkiw (2018) | Duolingo                                 | Vocabulary  | U.S.                     | Spanish          | No significant difference between control group and experimental group were found.                   |
| Castañeda and Cho (2016)           | Conjugation Nation                       | Grammar   | U.S.                     | Spanish          | The group improved skills after traditional teaching and app use.                                    |
| Ebadi and Ghuchi (2018)            | Memrise                                  | Vocabulary  | Iran                     | English          | The group using the app performed better than the group using traditional teaching methods.          |
| Fouz-González (2020)               | English File Pronunciation Spaceteam EFL | Pronunciation                                     | Spain                    | English          | No significant difference was found between app group and traditional teaching group.                |
| Grimshaw and Cardoso (2018)        | Spaceteam EFL                            | Oral Fluency                                      | Canada (French speaking) | English          | No significant difference existed between the groups in either rate of speech or fluency.            |
| Ou-Yang and Wu (2017)              | My Eva Mobile                            | Vocabulary  | Taiwan                   | English          | English majors performed better than non-English majors.   |
| Shih (2017)                        | LINE                                     | Listening Comprehension and Reading Comprehension | Taiwan                   | English          | The group using the app performed better than those using a language laboratory.                     |

Pre- and post-tests allowed participants to act as their own control groups. Results showed statistically significant improvements were made for both present indicative and present subjunctive conjugations. However, the authors were unable to determine whether the app contributed to increased knowledge gains compared to the traditional teaching methods.

Memrise was the app studied by Ebadi and Ghuchi (2018) regarding English vocabulary learning. For four weeks, 40 college students in Iran were divided into a control group that received traditional face-to-face instruction and an experimental group that received the same instruction as the control group plus the use of Memrise outside of class. The participants completed a 40-question pre- and post-test on vocabulary which showed that students using Memrise scored statistically significantly higher than students who just received face-to-face instruction ( $p < .05$ ). Interviews completed with a subsection of participants indicated they used the app between two and three and a half hours per week during their free time.

Fouz-González (2020) studied the English File Pronunciation App's ability to improve the pronunciation skills of 52 Spanish learners studying English at a university in Spain. At the beginning of the study, both an experimental and a control group were formed. Pre- and post- tests were given to both groups in three areas, perception, identification, and production two weeks later. Unfortunately, apart from the control group not using the app for the first two weeks, no information was given regarding whether the control group received any teaching between the pre- and post-test. To provide the same educational opportunities for all, the researchers changed the control group into an experimental group after two weeks, but only for the perception tasks. Results were mixed. No significant difference was found in the discrimination tasks, although there was improvement for both groups. The experimental group, however, improved significantly more than the control group in identification. Note that once the control group became the experimental group, their scores in identification improved by a factor of three compared to when they were not using the app. In the production tasks, the experimental group showed a significant difference in post-test scores for sentence reading and picture description, but not for imitation, although within the imitation task, a significant improvement was seen for the phoneme/ae/. The researchers did not complete a post-test on the new experimental group "to avoid imposing excessive demands on participants" so no results are available.

Grimshaw and Cardoso (2018) asked whether the Spaceteam ESL app could improve oral fluency as measured by syllables produced per minute and judge's ratings. Twenty French speaking Canadian university students were divided into two groups, a control group that was taught via

656 traditional face-to-face methods, and an experimental group which  
657 started each class with 15 minutes of app use over six weeks. Participants  
658 using the app worked in teams and were required to provide clear  
659 information to each other during the game. Pre- and post-tests showed  
660 no significant difference between the groups in either rate of speech or  
661 fluency as measured through judges' ratings, although the experimental  
662 group did score statistically significantly higher on their delayed post-test  
663 compared to their pre-test ( $p < .05$ ).

664 Shih (2017) studied the effectiveness of the LINE app on English  
665 listening comprehension and reading comprehension with 72 university  
666 juniors in Taiwan. The control group received traditional English instruc-  
667 tion and participated in a language lab over the course of ten weeks  
668 while the experimental group received direct instruction and used the  
669 mobile app. Pre- and post-tests showed that the experimental group  
670 performed significantly better than the control group.

671 Finally, Ou-Yang and Wu's study (2017) was different from the others  
672 described above in that it did not compare groups with and without an  
673 app. Instead, it used the same app, MyEva Mobile, to compare different  
674 groups on their vocabulary learning. The two groups consisted of 55  
675 English majors and 53 non-English majors at a university in Taiwan.  
676 Results showed that while both groups improved after using the app  
677 for two weeks, the English majors improved significantly more than the  
678 non-English majors.

## 680 Discussion

681  
682 The goal of this study was to address our five research questions. A  
683 discussion of each research question follows in addition to a more  
684 general discussion of efficacy, methodology, limitations, and  
685 recommendations.

686 The first research question asked how many scientifically reputable  
687 studies exist on FLL app efficacy. The most important finding linking  
688 the current study to the efficacy studies of Sung et al. (2015) and Burston  
689 (2015) was discovering of the extreme paucity of rigorous studies com-  
690 pleted on the efficacy of mobile apps for both FLL and MALL in general.  
691 Burston (2015) laments that "statistically reliable measures of learning  
692 outcomes are few and far between" in MALL (p. 16) while Sung et al.  
693 (2015) states that there is a "dearth of review research into the effec-  
694 tiveness of MALL and mobile-device-assisted teaching" (p. 8). As FLL  
695 apps are a subcategory of MALL, these statements are even more rep-  
696 resentative of app efficacy. The fact that only eight studies met this  
697 review's criteria has important implications for what we know about  
698 FLL app efficacy which will be explored later in the discussion.

699 The second question investigated what general trends exist in the  
700 research. Overall, the publication of FLL app efficacy studies began with  
701 only one study in 2016 with no real growth to 2020. Studies were spread  
702 internationally with seven countries accounting for the publication of  
703 the eight included articles. The amount of time the studies covered  
704 ranged from two weeks to four semesters, thereby precluding any infor-  
705 mation about using apps for long-term language learning.

706 While Sung et al.'s review (2015) found that elementary school-age  
707 students were the most commonly studied group in MALL, the current  
708 study of FLL apps found college students to be more frequently repre-  
709 sented, which was closer to Burston's (2015) MALL review which found  
710 that the most common age group of participants was 25-49 years old.  
711 Burston and Sung et al.'s studies were both in the field of MALL and  
712 represented the same approximate timeframe. However, Sung et al.  
713 included 44 studies which comprised of peer-review journal articles,  
714 unpublished conference presentations, and doctoral dissertations, whereas  
715 Burston's study used the MALL implementation database as a source,  
716 also allowing for publications such as conference proceedings, but with  
717 stringent inclusion criteria. The current study differed from these in  
718 that it limited the area of inquiry from MALL to FLL mobile apps and  
719 required studies to be published in peer-review journals. Given the small  
720 number of studies meeting the criteria, it may be more meaningful to  
721 consider the field of MALL when gathering information on specific age  
722 groups that have been studied. Note that Sung et al. had college students  
723 listed as the second most common group of participants studied after  
724 elementary students with percentages of 26.7 and 40 respectively.

725 The third research question focused on which languages were taught  
726 and to what linguistic audiences. Native speakers of five different lan-  
727 guages, French, Persian, English, Spanish and Chinese, were participants  
728 in the studies. Concurring with Sung et al. (2015) findings that English  
729 as a target language was the most studied, five of the eight studies  
730 focused on English learning with the other languages being Spanish and  
731 French. This is representative of English being the most studied foreign  
732 language worldwide (Serrato & Rodriguez, 2020). However, it is con-  
733 cerning that even for English, only five studies met our criteria. The  
734 problem lies herein; if it takes several studies to show that a particular  
735 mobile app is effective in teaching certain areas of language, this may  
736 only apply to learners who are at a certain level. Add to this the fact  
737 that each app needs to be tested for efficacy in several areas of language  
738 such as pronunciation, grammar, and listening comprehension. Next,  
739 imagine that research shows that a company has an app that is effective  
740 in teaching English to native Spanish speakers. It will not automatically  
741 follow that another version of this app, such as one teaching English

742 to Japanese speakers or teaching Russian to English speakers will be as  
743 effective. The amount of research required to evaluate the efficacy of a  
744 wide variety of mobile apps in FLL is nearly non-existent compared to  
745 what is needed.

746 Although we could have considered changing our inclusionary criteria  
747 to allow for the inclusion of a wider selection of studies regarding FLL  
748 app efficacy, we decided to keep our original inclusionary criteria for  
749 two reasons. First, Burston maintains high demands for rigorous  
750 research, although we did omit two criteria deemed unnecessary.  
751 Second, the nature of systematic reviews requires selection criteria to  
752 be finalized and possibly published before the study is carried out,  
753 ensuring that methods are not changed after the fact. Moher et al.  
754 (2015) stated that “Ideally, systematic reviews are based on pre-defined  
755 eligibility criteria and conducted according to a pre-defined method-  
756 ological approach” (p. 1).

757 Regarding the fourth research question of what specific aspects of  
758 linguistic knowledge and skills were measured, the current study con-  
759 cluded, similar to Dehghanzadeh et al. (2019), Heil et al. (2016), Sung  
760 et al. (2015), Burston (2015), and Duman et al. (2015), that vocabulary  
761 was the most represented language area studied in FLL, both by apps  
762 and overall MALL. Out of eight studies meeting the criteria for inclu-  
763 sion, three had vocabulary as its only measure, two of which measured  
764 improvement in English vocabulary for participants speaking Arabic  
765 (Ebadi & Ghuchi, 2018) and Chinese (Ou-Yang & Wu, 2017), and one  
766 which examined the acquisition of Spanish vocabulary for participants  
767 in the U.S. While most FLL mobile apps and MALL programs were not  
768 limited to teaching vocabulary, results showed that this area was the  
769 most often measured. Although vocabulary is certainly necessary to  
770 learn a language, it is but one area of language. Perhaps the reason for  
771 this focus on vocabulary learning is the relatively easy possibility of  
772 measuring the acquisition of a vocabulary item in a binary fashion  
773 within an app, such as measuring whether a word was correctly chosen  
774 from a group or produced from memory. It is certainly more difficult  
775 to develop a study that measures the correctness of pronunciation, for  
776 example, as measures are somewhat subjective and progress is made in  
777 increments, not in a binary fashion of being correct or incorrect.  
778 Grammar is also a relatively difficult area to measure, given the high  
779 number of grammatical components that exist and also interact with  
780 each other. Although a specific exercise such as forming the past tense  
781 of a given verb could be assessed with a binary scoring system, getting  
782 an overall measure of grammatical skill is a complex task. Similarly,  
783 skills such as oral fluency, listening comprehension, and reading ability  
784 are likely to be considered as more daunting to measure. It is also

785 possible that the explicit nature of vocabulary learning provides app  
786 users with the feeling of greater face validity and perceived efficacy.

787 Although few apps are limited to vocabulary learning, Heil et al.  
788 (2016) reported in their study of a selection of 50 mobile apps available  
789 in Google Play or the Apple iTunes App Store, that 42 out of the 50  
790 emphasized teaching vocabulary items as isolated units without context.  
791 This is tempered by the fact that 23 of the 50 were also found to  
792 emphasize vocabulary in context, meaning that only 19 of the 50 were  
793 limited to isolated vocabulary items. The researchers asserted that this  
794 was still problematic as language is a communicative tool that requires  
795 learners to know how to meaningfully use the words they have learned.

796 One possible way of addressing this in future research is to find  
797 measures of full language competency including various areas of language  
798 knowledge and skills which could be used in any FLL app, thereby  
799 providing a basis for comparison. Suggested guidelines for foreign lan-  
800 guage assessment have been proposed by groups such as the European  
801 Association for Language Testing and Assessment (EALTA, 2006), the  
802 Japanese Language Testing Association (JLTA, 2002), and the International  
803 Language Testing Association (ILTA, 2007). Evaluating specific guidelines  
804 in terms of mobile apps is beyond the realm of this discussion but we  
805 encourage the FLL community to further explore this area.

806 The fifth and final research question asked what the efficacy outcomes  
807 of these studies were. Considering the very small number of studies  
808 examined in this review, an overall analysis of the studies' outcomes  
809 does not lead to a simple interpretation. First, we examined the efficacy  
810 of the two Duolingo studies to try to make a direct comparison. However,  
811 each study focused on different areas of language (listening, reading,  
812 and writing versus vocabulary), tested different age groups (college  
813 students versus third and fourth graders) with different native languages  
814 (Spanish versus English) while attempting to learn different languages  
815 (French versus Spanish). The studies also greatly differed in research  
816 design with one analyzing whether self-regulation and scaffolding train-  
817 ing used within the app was more effective compared to the app without  
818 self-regulation training. In contrast, in the other study, the control group  
819 received traditional teaching while the experimental group only used  
820 the app for foreign language learning. García Botero et al. (2019) con-  
821 cluded that app use in addition to traditional teaching was not useful  
822 in terms of learning measures unless combined with self-regulation and  
823 scaffolding activities. However, Rachels and Rockinson-Szapkiw (2018)  
824 study indicated that Duolingo was equally effective as a traditional  
825 curriculum. Caution must be used when interpreting the results. Their  
826 study limits its evaluation of Duolingo to vocabulary although the app  
827 teaches other aspects of language too. This limits the study's claims to

828 one of many language areas. Note that the face-to-face curriculum of  
829 the control group was based on the content of Duolingo, thereby asking  
830 the question of whether the Duolingo vocabulary content could be taught  
831 better through traditional instruction or through the app. In this case,  
832 the answer appears to be the latter. Furthermore, the study uses a con-  
833 venience sample as opposed to a randomized sample, making it impos-  
834 sible to say that the results are generalizable to other groups. Strong  
835 differences in study design also made it difficult to compare the only  
836 two studies that exist for the same app.

837 It should be noted that while the ACT Framework promoted mea-  
838 surement of efficacy in the skills the instructional tool was designed to  
839 improve (Mattern, 2019), commercial apps were not designed by the  
840 researchers and generally do not come with information about what  
841 specific areas of language-learning they target. This means that the  
842 failure of an app to improve learning outcomes in a specific area such  
843 as vocabulary in no way means that the app is not useful in other areas  
844 of language-learning.

845 The range of methodologies used continued to grow with the number  
846 of studies examined. While the LINE, Memrise, and Spaceteam EFL  
847 apps were relatively similar in that they compared app use to traditional  
848 teaching practices, others had very different methods and goals. For  
849 example, students using Conjugation Nation all received face-to face  
850 instruction along with app use, and acted as their own controls, but  
851 without the possibility of attributing success to a particular part of the  
852 instruction. The study of MyEva Mobile differed from all the others by  
853 comparing different groups using the same app, in this case English  
854 majors and non-English majors, to see if one was more successful than  
855 the other, thereby telling us relatively little about whether the app was  
856 successful relative to traditional teaching. Furthermore, it is possible  
857 that the English-majors had more English-language input before and  
858 during the experiment, leading to this group having more exposure to  
859 the vocabulary items which in turn led to higher vocabulary scores  
860 (Peters, 2018). It should be noted that in several studies, the researchers  
861 gave participants instructions on how to use the apps, five out of eight  
862 studies were used only during class and some were used with students  
863 working in groups. This is certainly a different context compared to  
864 that of an average consumer who may download an app on a tablet  
865 device and then use it alone at times of their choice.

866 The overall results of efficacy studies will be important because as  
867 more research is carried out in this area, consumers will be well served  
868 by scientific information about the efficacy of FLL apps. Studying a  
869 foreign language is time consuming, therefore, app users may be inter-  
870 ested in knowing which apps have been shown to be effective. Also,

871 very specific efficacy information could be useful for individuals who  
872 want to focus on a specific area of language such as vocabulary instead  
873 of working on more global areas. For instance, a focus on listening as  
874 a first skill is recommended by Gangaiamaran and Pasupathi (2017)  
875 who make the case that this is the first skill that infants use during the  
876 acquisition of their native language. They also correctly note that lis-  
877 tening to a foreign language after achieving a certain level of proficiency  
878 is even more difficult due to the natural instinct to focus attention not  
879 only on the sounds of the language, but also on the meaning they carry.  
880 This leads to the interesting question of whether in some specialized  
881 cases, apps focusing on a single skill such as listening may be pedagog-  
882 ically desirable, leading to ensuing app use combining the teaching of  
883 several linguistic skills. As the acquisition of several linguistic skills is  
884 necessary to create fluency in a language, we recommend interdisciplin-  
885 ary collaborations between app designers, linguists, and academic  
886 researchers to develop more efficacious FLL apps and improve FLL app  
887 quality for consumers.

888 One more group that will be informed by these results are researchers  
889 in MALL. Hopefully, the realization of the extremely limited number  
890 of studies that met the criteria of this systematic review will spur the  
891 field to design studies that are scientifically rigorous and to publish  
892 them in peer-reviewed journals. Database searches turn out very high  
893 numbers of publications on FLL apps, but the majority of them, although  
894 many are likely to be of high quality, are published in other venues,  
895 lacking the gold standard that these journals have to offer.

## 897 **Limitations**

898 Limitations of this systematic review must be acknowledged. First, no  
899 review of this type can claim to be completely exhaustive since search  
900 terms and databases must be limited by some constraints, some of which  
901 are due to be imperfect. We also acknowledge that several efficacy  
902 studies exist that are not listed here. For example, a 2012 final white  
903 paper on the efficacy of Duolingo by Vesselinov and Grego (2016)  
904 reported findings from an experiment of 196 adults studying Spanish  
905 for eight weeks. Results claimed that a person with no knowledge of  
906 Spanish would be able to reach a degree of knowledge of the language  
907 equivalent to a college semester after an average of 34 hours of study  
908 with the app. However, this article was not published in a peer-reviewed  
909 academic journal. Second, as noted by Haidich (2010), a meta-analysis  
910 requires that studies use the same designs, have similar participants,  
911 and have the same outcomes. In our case, since studies had different  
912 outcomes, used different participants, different mobile apps, a  
913

914 meta-analysis would not be possible. Therefore, we completed a system-  
915 atic review and were unable to compare effect sizes between studies.

## 916 917 **Recommendations**

918  
919 Based on this study, there are multiple recommendations. The *first rec-*  
920 *ommendation has two parts*. First, FLL app designers should be encour-  
921 aged to build apps that focus on communicative ability as a whole. The  
922 strong focus on vocabulary teaching and measurement has been shown  
923 repeatedly. Out of the eight studies discussed here, three of them focused  
924 on vocabulary teaching. Dehghanzadeh et al. (2019) research, which  
925 examined 22 FLL games, found that 17 were focused on vocabulary.  
926 Vocabulary was also the most commonly studied area in Sung et al.  
927 (2015) MALL efficacy meta-analysis and Burston's (2015) review of  
928 MALL learning outcomes. Heil et al. (2016) reviewed the 50 most pop-  
929 ular FLL apps, to find that 42 taught vocabulary in isolation. This review  
930 also showed vocabulary to be the most tested area. This suggests that  
931 many commercial FLL apps focus on vocabulary, often independently  
932 of other language areas. The thesis that vocabulary should not be taught  
933 in isolation was put forward not only in pedagogical terms but in psy-  
934 cholinguistic terms by Bolgün and McCaw (2019) who evaluated lan-  
935 guage technology in light of our understanding of the neuroscience of  
936 memory, particularly the declarative and procedural memory systems.  
937 They argue that "language technology caters predominantly to the declar-  
938 ative memory system" involved in explicit memorized knowledge, whereas  
939 it should "cater to the procedural memory" system involved in implicit  
940 grammatical processes like grammar and sentence structure (Bolgün &  
941 McCaw, 2019). The idea is that explicit knowledge of vocabulary does  
942 not equate to an implicit understanding of its use conversationally.

943 The second part of this recommendation is closely related to the first.  
944 Tools measuring the overall communicative competence of foreign lan-  
945 guage learners must be developed and used both within and across  
946 studies to compare app efficacy in meaningful ways. Just as standardized  
947 IQ tests provide general measures of intelligence, standardized language  
948 tests would allow researchers to more easily compare results across  
949 future studies. It is important to note, however, that just because studies  
950 may choose to examine a single aspect of language, such as vocabulary,  
951 does not mean the app is limited to that aspect. Developing measures  
952 of overall communicative efficacy would be an important step forward  
953 in expanding the field of FLL efficacy.

954 *Our next recommendation* is for FLL researchers to measure learning  
955 outcomes in a scientific, rigorous manner. Perhaps the most important  
956 finding from this systematic review is the realization that despite the

957 large number of returned results when searching for studies about FLL  
958 apps, only a tiny percentage of these measured efficacy and met the  
959 criteria of this study, which were more relaxed than those put forward  
960 by Burston (2015). This means that the number of high-quality scientific  
961 studies measuring the effect that commercial FLL apps have on learning  
962 outcomes is shockingly small. As stated above, this needs to change for  
963 the good of FLL students, app consumers, app developers, and FLL  
964 instructors. Given the large number of commercially available FLL apps,  
965 there is much room for further exploration into their efficacy. While it  
966 is certainly valuable to examine these apps in other contexts such as  
967 exploring user enjoyment, motivation, and usage trends, it is vital to  
968 the field that the efficacy of MALL tools, in this case apps, is measured  
969 in systematic ways that allow studies to be compared.

970 *We also recommend* that studies be developed that compare several  
971 different experimental groups, each using a different FLL app only, to  
972 control groups using clearly defined traditional instruction. Measures  
973 would be taken globally or by testing different areas of language includ-  
974 ing pronunciation, listening comprehension, vocabulary, and grammar.  
975 This would be of use to consumers who have decided to use an app  
976 but need assistance in deciding which one(s) to choose. Also useful,  
977 particularly for FLL educators, would be studies where different apps  
978 are combined with identical traditional methods. Furthermore, compar-  
979 isons of app efficacy when instructions are given versus when the user  
980 is left to their own devices would be helpful. This would allow FLL  
981 instructors using apps as part of their teaching to know whether instruc-  
982 tion in app use is to be recommended. In summary, there is an oppor-  
983 tunity in the field to continue researching the efficacy of mobile  
984 applications so educators, students, and consumers can be better informed  
985 when making decisions about applications designed to teach foreign  
986 **Q3** languages.

### 987 **Disclosure statement**

988 **Q4** No potential conflict of interest was reported by the authors.

### 989 **ORCID**

990 Chrystal Sapphire Dragonflame  <http://orcid.org/0000-0001-6688-5807>

991 Amanda A. Olsen  <http://orcid.org/0000-0002-7707-7271>

### 992 **References**

993 Babel. (2020, December 21). *About us*. <https://about.babbel.com/en/about-us/>

**Q6**

- 1000 Bachman, L. F., & Palmer, A. S. (1996). *Language testing in practice: Designing and*  
1001 *developing useful language tests*. Oxford University Press.
- 1002 Bolgün, M. A., & McCaw, T. (2019). Toward a neuroscience-informed evaluation of  
1003 language technology. *Computer Assisted Language Learning*, 32(3), 294–321. <https://doi.org/10.1080/09588221.2018.1516675>
- 1004 Burston, J. (2015). Twenty years of MALL project implementation: A meta-analysis of  
1005 learning outcomes. *ReCALL*, 27(1), 4–20. <https://doi.org/10.1017/S0958344014000159>
- 1006 Castaneda, D. A., & Cho, M. H. (2016). Use of a game-like application on a mobile  
1007 device to improve accuracy in conjugating Spanish verbs. *Computer Assisted Language*  
1008 *Learning*, 29(7), 1195–1204. <https://doi.org/10.1080/09588221.2016.1197950>
- 1009 Chan, J. (2019, January 31). Top language apps worldwide for 2018 by downloads.  
1010 *Sensortower*. <https://sensortower.com/blog/top-language-apps-2018-ww>
- 1011 Clearinghouse, W. W. (2012). What works clearinghouse. *Internet site*. <http://ies.ed.gov/ncee/wwc>.
- 1012 Dehghanzadeh, H., Fardanesh, H., Hatami, J., Talaei, E., & Noroozi, O. (2019). Using  
1013 gamification to support learning English as a second language: A systematic review.  
1014 *Computer Assisted Language Learning*, 1–24. [https://doi.org/10.1080/09588221.2019.1](https://doi.org/10.1080/09588221.2019.1648298)  
1015 [648298](https://doi.org/10.1080/09588221.2019.1648298) Q7
- 1016 Doabler, C. T., Clarke, B., Kosty, D., Turtura, J. E., Firestone, A. R., Smolkowski, K.,  
1017 ... Maddox, S. A. (2019). Efficacy of a first-grade mathematics intervention on  
1018 measurement and data analysis. *Exceptional Children*, 86(1), 77–94. <https://doi.org/10.1177/0014402919857993>
- 1019 Duman, G., Orhon, G., & Gedik, N. (2015). Research trends in mobile assisted language  
1020 learning from 2000 to 2012. *ReCALL*, 27(2), 197–216. [https://doi.org/10.1017/](https://doi.org/10.1017/S0958344014000287)  
1021 [S0958344014000287](https://doi.org/10.1017/S0958344014000287)
- 1022 Duolingo. (2020, December 21). *Free language courses for English speakers*. <https://www.duolingo.com/courses>
- 1023 Dynarski, M. (2015). Using research to improve education under the Every Student  
1024 Succeeds Act. *Evidence Speaks Reports*, 1(8), 1–5.
- 1025 Ebadi, S., & Ghuchi, K. D. (2018). Investigating the effects of blended learning approach  
1026 on vocabulary enhancement from EFL learners' perspectives. *Journal on English*  
1027 *Language Teaching*, 8(2), 57–68. <https://doi.org/10.26634/jelt.8.2.13981>
- 1028 Fouz-González, J. (2020). Using apps for pronunciation training: An empirical evaluation  
1029 of the English File Pronunciation app. *Language Learning & Technology*, 24(1),  
1030 62–85. <http://hdl.handle.net/10125/44709>
- 1031 Gangaianaran, R., & Pasupathi, M. (2017). Review on use of mobile apps for language  
1032 learning. *International Journal of Applied Engineering Research*, 12(21), 11242–11251.  
1033 <https://doi.org/10.1254/44709>
- 1034 García Botero, G., Botero Restrepo, M. A., Zhu, C., & Questier, F. (2019). Complementing  
1035 in-class language learning with voluntary out-of-class MALL. Does training in  
1036 self-regulation and scaffolding make a difference? *Computer Assisted Language*  
1037 *Learning*, 1–27. <https://doi.org/10.1080/09588221.2019.1650780> Q8
- 1038 Goldman, S. R., Greenleaf, C., Yukhymenko-Lescroart, M., Brown, W., Ko, M. L. M.,  
1039 Emig, J. M., ... Britt, M. A. (2019). Explanatory modeling in science through  
1040 text-based investigation: Testing the efficacy of the Project READI intervention ap-  
1041 proach. *American Educational Research Journal*, 56(4), 1148–1216. <https://doi.org/10.3102/0002831219831041>
- 1042 Grimshaw, J., & Cardoso, W. (2018). Activate space rats! Fluency development in a  
mobile game-assisted environment. *Language Learning & Technology*, 22(3), 159–175.  
<https://core.ac.uk/download/pdf/211326392.pdf>

- 1043 Haidich, A. B. (2010). Meta-analysis in medical research. *Hippokratia*, 14(Suppl 1),  
 1044 29–37. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3049418/>
- 1045 Heil, C. R., Wu, J. S., Lee, J. J., & Schmidt, T. (2016). A review of mobile language  
 1046 learning applications: Trends, challenges, and opportunities. *The EuroCALL Review*,  
 1047 24(2), 32–50. <https://doi.org/10.4995/eurocall.2016.6402>
- 1048 Koedinger, K. R., Corbett, A. T., & Perfetti, C. (2012). The knowledge-learning-instruction  
 1049 framework: Bridging the science-practice chasm to enhance robust student  
 1050 learning. *Cognitive Science*, 36(5), 757–798.
- 1051 Mattern, K. (2019). *ACT's efficacy framework: The intersection of learning, measurement,  
 1052 and navigation. issue brief*. ACT, Inc.
- 1053 Moher, D., Shamseer, L., Clarke, M., Ghersi, D., Liberati, A., Petticrew, M., Shekelle,  
 1054 P., & Stewart, L. A. (2015). Preferred reporting items for systematic review and  
 1055 meta-analysis protocols (PRISMA-P) 2015 statement. *Systematic Reviews*, 4(1), 1–9.  
 1056 <https://doi.org/10.1186/2046-4053-4-1>
- 1057 Q9 Olsen, A., Dragonflame, C., & Tommerdahl, J. (Under Review). *A systematic review*.
- 1058 Ou-Yang, F. C., & Wu, W. C. V. (2017). Using mixed-modality vocabulary learning on  
 1059 mobile devices: Design and evaluation. *Journal of Educational Computing Research*,  
 1060 54(8), 1043–1069. <https://doi.org/10.1177/0735633116648170>
- 1061 Pandey, M., Litoriya, R., & Pandey, P. (2019). Perception-based classification of mobile  
 1062 apps: A critical review. In A. K. Luhach, K. B. G. Hawari, I. C. Mihai, P. Hsiung,  
 1063 & R. B. Mishra (Eds.), *Smart computational strategies: Theoretical and practical aspects*  
 1064 (pp. 121–133). Springer. [https://doi.org/10.1007/978-981-13-6295-8\\_11](https://doi.org/10.1007/978-981-13-6295-8_11)
- 1065 Peters, E. (2018). The effect of out-of-class exposure to English language media on  
 1066 learners' vocabulary knowledge. *Itl - International Journal of Applied Linguistics*,  
 1067 169(1), 142–168. <https://doi.org/10.1075/itl.00010.pet>
- 1068 Purpura, J. E. (2013). Assessing grammar. *The Companion to Language Assessment*, 1,  
 1069 100–124. <https://doi.org/10.1002/9781118411360.wbcla147>
- 1070 Rachels, J. R., & Rockinson-Szapkiw, A. J. (2018). The effects of a mobile gamification  
 1071 app on elementary students' Spanish achievement and self-efficacy. *Computer Assisted  
 1072 Language Learning*, 31(1-2), 72–89. <https://doi.org/10.1080/09588221.2017.1382536>
- 1073 Rovio Entertainment Corporation. (2020, March 2). *About - Rovio*. <https://www.rovio.com/about/>
- 1074 Serrato, D. I., & Rodriguez, B. C. P. (2020). Academic e-tandems as a strategy for  
 1075 English language learning in a Mexican university. *Open Praxis*, 12(3), 417–424.  
 1076 <https://doi.org/10.5944/openpraxis.12.3.1099>
- 1077 Shamseer, L., Moher, D., Clarke, M., Ghersi, D., Liberati, A., Petticrew, M., Shekelle,  
 1078 P., Stewart, L. A., & Group, P.-P. (2015). Preferred reporting items for systematic  
 1079 review and meta-analysis protocols (PRISMA-P) 2015: Elaboration and explanation.  
 1080 *BMJ*, 349 <https://doi.org/10.1136/bmj.g7647>
- 1081 Q10 Shih, R. C. (2017). The effect of English for Specific Purposes (ESP) learning-language  
 1082 lab versus mobile-assisted learning. *International Journal of Distance Education  
 1083 Technologies (IJDET)*, 15(3), 15–30. <https://doi.org/10.4018/IJDET.2017070102>
- 1084 Sriganesh, K., Shanthanna, H., & Busse, J. W. (2016). A brief overview of systematic  
 1085 reviews and meta-analyses. *Indian Journal of Anaesthesia*, 60(9), 689–694. <https://doi.org/10.4103/0019-5049.190628>
- Sung, Y. T., Chang, K. E., & Yang, J. M. (2015). How effective are mobile devices for  
 language learning? A meta-analysis. *Educational Research Review*, 16, 68–84. <https://doi.org/10.1016/j.edurev.2015.09.001>
- Toste, J. R., Capin, P., Williams, K. J., Cho, E., & Vaughn, S. (2019). Replication of an  
 experimental study investigating the efficacy of a multisyllabic word reading inter-

- 1086           vention with and without motivational beliefs training for struggling readers. *Journal*  
1087           *of Learning Disabilities*, 52(1), 45–58. <https://doi.org/10.1177/0022219418775114>  
1088 Toto, G. A., & Limone, P. (2019). Contemporary trends in studies on mobile learning  
1089           of foreign languages: A meta-analysis. *International Journal of Engineering Education*,  
1090           1(2), 85–90. <https://doi.org/10.14710/ijee.1.2.85-90>  
1091 Vaughn, S., Martinez, L. R., Williams, K. J., Miciak, J., Fall, A. M., & Roberts, G.  
1092           (2019). Efficacy of a high school extensive reading intervention for English learners  
1093           with reading difficulties. *Journal of Educational Psychology*, 111(3), 373. [https://doi.](https://doi.org/10.1037/edu0000289)  
1094           [org/10.1037/edu0000289](https://doi.org/10.1037/edu0000289)  
1095 Vesselinov, R., & Grego, J. (2016). *The Babel efficacy study: Final report* [White paper].  
1096           Babel. [https://press.babel.com/fr\\_CA/releases/downloads/Babel-Efficacy-Study.pdf](https://press.babel.com/fr_CA/releases/downloads/Babel-Efficacy-Study.pdf)  
1097 Wanzek, J., Vaughn, S., Roberts, G., & Fletcher, J. M. (2011). Efficacy of a reading  
1098           intervention for middle school students Identified with Learning Disabilities.  
1099           *Exceptional Children*, 78(1), 73–87. <https://doi.org/10.1177/001440291107800105>  
1100 Yu, Z. (2019). A systematic review on mobile technology-assisted English learning.  
1101           *International Journal of e-Collaboration*, 15(4), 71–88. [https://doi.org/10.4018/](https://doi.org/10.4018/IJeC.2019100105)  
1102           [IJeC.2019100105](https://doi.org/10.4018/IJeC.2019100105)  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128