

1 **Analyse psychométrique d'outils d'évaluation mathématique utilisés auprès**
2 **des enfants francophones**

3

4 Psychometric Analysis of Mathematic Assessment Tools Used with French-
5 speaker Children

6

7 **Anne Lafay, Ph.D.**

8 Orthophoniste

9 Chercheuse postdoctorale

10 Email : lafay_anne@yahoo.fr

11

12 **Julie Cattini**

13 Orthophoniste

14 Email : juliecattini@gmail.com

15

16

17 **Notes des auteurs** : Une des auteures (Anne Lafay) est une des co-auteurs d'un outil de
18 tests appartenant à la recension des outils d'évaluation mathématique et qui fait l'objet de
19 l'analyse.

20

21 **Abrégé**

22 Si nous nous référons au *Manuel diagnostique et statistique des troubles mentaux* (5^e éd.;
23 DSM-5; American Psychiatric Association, 2013 [version anglaise], 2016 [version
24 française]), l'évaluation d'un enfant en difficulté mathématique doit comporter une évaluation
25 objective. Cette évaluation vient aider le clinicien à déterminer si les compétences scolaires de
26 l'enfant sont nettement en-dessous du niveau escompté pour l'âge chronologique. Jusqu'à
27 présent, aucune étude ne s'est intéressée à évaluer les qualités psychométriques des tests
28 disponibles en français pour évaluer les capacités mathématiques des enfants francophones.
29 Pourtant, les professionnels sont amenés à faire un choix éclairé sur le ou les test(s) qu'ils
30 utiliseront avec leurs patients. La présente étude vise, d'une part, à mettre à jour la recension
31 des outils disponibles en français pour l'évaluation mathématique qui avait été établie par
32 Lafay, St-Pierre et Macoir (2014) et, d'autre part, à analyser leurs qualités psychométriques.
33 Les résultats obtenus démontrent que, bien que plusieurs outils soient disponibles, peu d'entre
34 eux répondent aux standards psychométriques. Cela remet donc en question la valeur
35 discriminante des outils disponibles. Ainsi, cette étude promeut l'utilisation d'une pratique
36 basée sur les données probantes pour aider les cliniciens à adopter une pratique réflexive lors
37 du choix des tests diagnostiques.

38
39 *Mots-clés* : évaluation, mathématique, trouble des apprentissages en mathématiques,
40 dyscalculie, enfant, francophone, psychométrie, validité, fidélité

41 L'évaluation d'un enfant en difficulté mathématique comporte un temps
42 d'investigation afin de mesurer objectivement les compétences scolaires. Jusqu'à présent,
43 aucune étude ne s'est intéressée à évaluer les qualités psychométriques des tests disponibles
44 en français pour évaluer les capacités mathématiques des enfants francophones. La présente
45 étude a pour objectif d'aider les orthophonistes à faire des choix éclairés dans la sélection des
46 outils d'évaluation des habiletés mathématiques. Ces renseignements pourront également leur
47 permettre d'être en mesure de mieux comprendre ou d'interpréter les résultats d'évaluations
48 complétées par d'autres professionnels.

49

50 **Diagnostic de trouble des apprentissages en mathématiques**

51 En fonction du milieu professionnel dans lequel il travaille, l'orthophoniste peut être
52 amené à travailler avec des enfants ayant un trouble de la communication ou qui sont aux
53 prises avec un problème de langage. Parmi ces enfants, plusieurs manifestent des difficultés
54 concomitantes en mathématiques. Par exemple, des difficultés en mathématiques ont été
55 observées chez les enfants sourds (pour une revue, voir Roux, 2014) et chez les enfants ayant
56 un trouble développemental du langage oral (Donlan, Cowan, Newton et Lloyd, 2007;
57 Durkin, Mok et Conti-Ramsden, 2013). Ajoutons que le trouble des apprentissages en
58 mathématiques – autrement appelé dyscalculie – est très fréquemment associé à la dyslexie.
59 En effet, entre 17% (Gross-Tsur, Manor et Shalev, 1996) et 43,3-65% (Barbaresi, Katusic,
60 Colligan, Weaver et Jacobsen, 2005) des enfants présentant un trouble des apprentissages en
61 mathématiques sont aussi dyslexiques selon les critères diagnostiques utilisés dans chacune
62 des études. Certains auteurs proposent que la capacité à traiter les nombres, qui sous-tend le
63 développement des habiletés mathématiques, est innée et disponible à tous, incluant aux
64 adultes sans culture et langage mathématique (Butterworth, Reeve, Reynolds et Lloyd, 2008;
65 Frank, Everett, Fedorenko et Gibson, 2008; Gordon, 2004), aux bébés (Antell et Keating,

66 1983; Starkey et Cooper, 1980; Wynn, 1992) et même aux animaux (Brannon, 2005).
67 D'autres proposent plutôt que les habiletés mathématiques se développent grâce à un système
68 numérique exact et lié spécifiquement au langage humain. À titre d'exemple, Carey (2001,
69 2004) attribue une importance primordiale au langage dans le développement du concept de
70 nombre. Selon cet auteur, même si les enfants perçoivent presque instantanément et de
71 manière quasi innée les très petites quantités (c.-à-d. *subitizing*), ce n'est que grâce à
72 l'acquisition des mots-nombres qu'ils deviennent capables d'associer une quantité précise à
73 un mot-nombre précis. En résumé, un fait demeure : le langage a une place prépondérante
74 dans le développement de la compréhension et dans l'application des concepts
75 mathématiques.

76 D'après la définition du *Manuel diagnostique et statistique des troubles mentaux* (5
77 éd.; DSM-5; American Psychiatric Association, 2013 [version anglaise], 2016 [version
78 française]), le trouble des apprentissages en mathématiques (1 à 10% des enfants d'âge
79 scolaire) est défini comme étant un déficit des apprentissages dans les sphères du sens du
80 nombre, du calcul et de la résolution de problèmes qui ne peut être expliqué par des troubles
81 d'ordre sensoriel, neurologique, psychiatrique ou environnemental. Le trouble des
82 apprentissages en mathématiques interfère fortement avec la scolarité des enfants en étant
83 atteints et avec les activités de la vie quotidienne impliquant des compétences numériques, et
84 ce, de manière durable et en dépit des interventions offertes. Les symptômes décrits dans le
85 DSM-5 sont, d'une part, des difficultés à maîtriser le sens des nombres, les données chiffrées
86 ou le calcul et/ou, d'autre part, des difficultés avec le raisonnement mathématique (p. ex. des
87 difficultés à appliquer des concepts ou des méthodes mathématiques pour résoudre les
88 problèmes). En s'appuyant sur cette définition, les cliniciens doivent alors inclure, dans
89 l'évaluation, une investigation minimale des habiletés mathématiques suivantes :
90 dénombrement (c.-à-d. action d'indiquer le nombre d'éléments d'une collection), lecture et

91 dictée de nombres (c.-à-d. action de passer d'un nombre écrit en code arabe [p. ex. 15] à un
92 nombre en code oral [p. ex. quinze], et inversement), calcul (p. ex. $2 + 4$, 3×12 , 150×3 ,
93 $1259 - 856$, etc.) et résolution de problèmes à énoncé verbal (p. ex. Marie-Ève a neuf jujubes
94 dans son sac. Elle en a trois de plus que son ami Marc-Antoine. Combien de jujubes Marc-
95 Antoine a-t-il?).

96 Le cadre théorique de l'approche cognitive (Butterworth, 2005; Noël et Rousselle,
97 2011; Von Aster et Shalev, 2007; Wilson et Dehaene, 2007) soutient par ailleurs que les
98 difficultés mathématiques mentionnées précédemment découlent d'une faiblesse au niveau du
99 *sens du nombre* (c.-à-d. du traitement des nombres présentés non symboliquement) et de
100 l'accès au *sens du nombre* via les nombres symboliques (c.-à-d. du traitement des nombres
101 présentés symboliquement). Cela suggère, par exemple, des difficultés à comparer des
102 nombres en format analogique (p. ex. ***) ou en format arabe (p. ex. 3), à identifier et à
103 estimer des quantités en format analogique, ou encore, à placer des nombres sur une ligne
104 numérique. En s'appuyant sur ces hypothèses théoriques, les cliniciens doivent alors inclure,
105 dans l'évaluation, une investigation minimale des habiletés du traitement des nombres
106 symboliques et non symboliques.

107 Selon le DSM-5, un enfant avec un trouble des apprentissages en mathématiques
108 présente un niveau mathématique significativement inférieur à ce qui est attendu pour l'âge,
109 tel qu'évalué par des tests standardisés de calcul et de raisonnement. Dans une démarche de
110 pratique basée sur les données probantes, on mesure ici toute l'importance d'utiliser les tests
111 les plus valides et de connaître, au préalable, leurs propriétés psychométriques afin de faire un
112 choix éclairé (Betz, Eickhoff et Sullivan, 2013; Gaul Bouchard, Fitzpatrick et Olds, 2009;
113 Leclercq et Veys, 2014; McCauley, 1989). Gaul Bouchard et al. (2009) ont par ailleurs
114 expliqué que « *l'Ordre des orthophonistes et audiologistes du Québec prévient leurs membres*
115 *que le fait de tirer des conclusions basées sur des tests non standardisés ou qui ne possèdent*

116 *pas des niveaux de qualités psychométriques appropriés va à l'encontre du code*
117 *déontologique de leur ordre professionnel* ». Betz et al. (2013) ont récemment investigué les
118 habitudes d'utilisation des tests de 364 cliniciens américains et ont rapporté la présence d'un
119 biais concernant la popularité d'un test. Ils ont en effet montré que la fréquence d'utilisation
120 des tests standardisés pour le diagnostic du trouble développemental du langage est
121 uniquement corrélée à la date de publication mais aucunement, malheureusement, à leurs
122 qualités psychométriques (p. ex. fidélité, validité, pouvoir discriminant). Ce résultat est tout à
123 fait surprenant et va à l'encontre des recommandations d'une pratique basée sur les données
124 probantes. De plus, il convient de reconnaître que les manuels de tests ne sont pas toujours
125 clairs et il est parfois difficile de s'y retrouver. Outre la lecture du manuel, il est donc
126 important de savoir ce que l'on cherche.

127

128 **Les qualités psychométriques d'un outil d'évaluation**

129 Un outil d'évaluation doit respecter plusieurs propriétés psychométriques pour être
130 considéré de bonne qualité : il doit être standardisé, valide, fidèle et posséder des données
131 normatives.

132

133 **Standardisation.** Un test est standardisé lorsque les conditions de passation et de
134 cotation ont été systématisées et uniformisées lors de l'étalonnage. De plus, le manuel est
135 suffisamment clair et précis pour permettre à un utilisateur ultérieur de reproduire de façon
136 identique les procédures et conditions dans sa pratique clinique. Ceci permet de limiter la
137 subjectivité, ainsi que les erreurs de mesure ou les biais d'interprétation. Vérifier les
138 caractéristiques de standardisation d'un test est primordial pour une utilisation optimale.

139

140 **Validité.** La validité d'un test réfère au degré avec lequel un test mesure vraiment ce
141 qu'il prétend mesurer. Plusieurs types de validité peuvent être investigués. Tout d'abord, la
142 validité de surface (ou validité d'apparence; Ivanova et Hallowell, 2013) est une mesure
143 subjective qui concerne la compréhension et l'acceptation du test par les utilisateurs (patients
144 et évaluateurs). Il s'agit de mesurer si l'évaluateur peut décrire l'objectif du test, s'il
145 comprend les consignes, s'il est capable d'utiliser le test et s'il juge que la présentation du test
146 est adéquate pour la tranche d'âge visée.

147 La validité de contenu (Gaul Bouchard et al., 2009), parfois appelée validité
148 théorique (Leclercq et Veys, 2014), réfère à la pertinence du contenu du test. On ne peut pas
149 affirmer qu'un test est valide pour toujours. La conception de l'outil et le choix des items qui
150 le composent doivent reposer sur les modèles théoriques récents de la fonction cognitive
151 évaluée.

152 La validité de critère (ou validité empirique; Ivanova et Hallowell, 2013) est la
153 capacité d'un test à évaluer adéquatement la performance par rapport à un critère de référence
154 (critère externe, indépendant). On distingue deux types de validation critériée : la validité
155 concomitante (autrement appelée validité concourante par Leclercq et Veys, 2014, ou encore,
156 validité concordante par Gaul Bouchard et al., 2009) et la validité prédictive. La validité
157 concomitante implique une comparaison, au même moment de mesure, entre le test et un
158 critère de référence externe (p. ex. un autre test standardisé mesurant le même construit
159 théorique). La validité prédictive implique une comparaison, en temps différé, entre le test et
160 un critère qui sert d'indicateur d'une performance future pour une tâche de nature similaire
161 que l'on cherche à prédire. La pertinence fonctionnelle de l'outil doit être attestée via une
162 concordance entre les scores observés à l'outil et le fonctionnement dans les activités de vie
163 quotidienne mettant en œuvre la fonction évaluée (p. ex. la note obtenue en mathématiques
164 lors des examens scolaires).

165 Enfin, la validité de construit réfère à la capacité d'un test ou d'une batterie de tests à
166 mesurer un construit théorique. Plusieurs analyses empiriques peuvent être rapportées dans les
167 manuels de test. Ces analyses devraient pouvoir s'expliquer en lien avec la théorie avancée
168 par le test. Ce n'est donc pas uniquement la valeur des analyses effectuées qu'il importe de
169 regarder, mais également son lien avec la définition du construit qu'il sous-tend. De ce fait, le
170 manuel d'un test doit, au préalable, définir les objectifs du test et des subtests de manière
171 claire et simple, en plus de définir le construit (c.-à-d. la définition conceptuelle, théorique ou
172 opérationnelle de ce qui est mesuré dans le test). La validité de construit est vérifiée par des
173 analyses portant sur la validité en lien avec les caractéristiques de l'individu, la validité
174 factorielle et la précision (Ivanova et Hallowell, 2013). La validité en lien avec les
175 caractéristiques de l'individu concerne le fait que lorsque le construit mesuré est
176 intrinsèquement relié à une ou plusieurs caractéristiques « évidentes » de l'individu, la mesure
177 de ce construit doit être sensible à cette relation (p. ex. sexe, âge, niveau socioéconomique,
178 pathologie, etc.). La validité factorielle est une mesure dans laquelle la structure théorique du
179 test correspond à la structure statistique observée. En d'autres mots, différents items ou sous-
180 tests, malgré des différences de contenu, de format ou de tâches, mesurent une dimension
181 commune qui influence la performance à tous les items ou sous-tests. Enfin, la précision ou le
182 pouvoir discriminant (également appelé pouvoir classificatoire; Ivanova et Hallowell, 2013)
183 d'un outil correspond à sa sensibilité et sa spécificité et doit garantir son pouvoir
184 diagnostique. La sensibilité est le pouvoir qu'un test possède pour repérer un enfant en
185 difficulté comme étant effectivement en difficulté (c.-à-d. un vrai positif). Selon Plante et
186 Vance (1994), un outil est reconnu comme étant sensible s'il permet de classer correctement
187 une forte proportion des personnes présentant des difficultés (80% à 95%). La spécificité est
188 le pouvoir qu'un test possède pour repérer une personne saine comme étant effectivement
189 saine (c.-à-d. un vrai négatif). Un outil est reconnu comme étant spécifique s'il permet de

190 classer correctement une forte proportion des personnes ne présentant pas de difficultés
191 (80% à 95%). En conclusion, la validité d'un instrument se détermine entre autres en évaluant
192 dans quelle mesure le test mesure réellement ce qu'il dit vouloir mesurer. Prendre
193 connaissance des différents éléments de validation d'un outil est une avenue indispensable
194 pour juger de sa pertinence dans un contexte particulier.

195

196 **Fidélité.** La fidélité d'un test porte sur son degré de cohérence, de précision et de
197 reproductibilité. Plusieurs types de fidélité peuvent également être investigués. Tout d'abord,
198 la cohérence interne (Ivanova et Hallowell, 2013) concerne le fait qu'un test psychologique
199 soit cohérent avec lui-même et que chacune de ses composantes réagisse de manière
200 cohérente à une même réponse. Il existe plusieurs analyses empiriques qui permettent
201 d'évaluer la cohérence interne. Tout d'abord, la cohérence interne peut être estimée par le
202 calcul du coefficient alpha de Cronbach. Il s'agit d'une valeur calculée qui s'étend entre 0 et
203 1. Plus la valeur alpha s'approche de 1, plus l'ensemble des éléments est homogène. Le seuil
204 minimal d'acceptabilité pour l'alpha de Cronbach est estimé à 0,70 (Nunnally, 1978). Il faut
205 toutefois noter qu'un alpha de Cronbach trop élevé peut être une indication de redondance. La
206 cohérence interne d'un test ou d'une batterie de tests peut également être évaluée à partir
207 d'une analyse de corrélations inter-items à l'intérieur d'une même épreuve, ou encore, d'une
208 analyse de corrélations inter-épreuves à l'intérieur d'une batterie de tests. Enfin, la cohérence
209 interne peut être estimée par une bissection des items (ou *split-half*), ce qui consiste à partager
210 aléatoirement un test en deux groupes d'items et à vérifier leur corrélation. Les balises
211 utilisées sont celles indiquées par Cohen (1988), à savoir qu'une corrélation autour de 0,10 est
212 faible, qu'une corrélation autour de 0,30 est moyenne et qu'une corrélation autour de 0,50 est
213 forte.

214 Ensuite, la fidélité peut être investiguée par une évaluation de la stabilité, qui consiste
215 à vérifier si le test donne des résultats relativement similaires (reproductibles) dans des
216 situations différentes et comparables. La fidélité temporelle (ou test-retest; Ivanova et
217 Hallowell, 2013) stipule que l'outil est en mesure de fournir des résultats comparables entre
218 deux passations à des temps différents, ce qui assure que les résultats obtenus ne sont pas
219 l'effet du hasard. La fidélité inter-juges (Ivanova et Hallowell, 2013) assure généralement que
220 les résultats obtenus par une personne sont le reflet de sa performance, indépendamment du
221 professionnel qui a administré et corrigé le test. Il importe donc que différents juges soient en
222 mesure d'évaluer les performances de la même manière. Ajoutons que lorsque deux versions
223 parallèles d'un même test existent, l'outil doit montrer que l'application de ces deux versions
224 aux mêmes personnes résultent en des scores équivalents (fidélité par versions parallèles;
225 Ivanova et Hallowell, 2013). En conclusion, la fidélité d'un instrument se détermine non
226 seulement en évaluant dans quelle mesure les items censés mesurer un même construit mènent
227 à des résultats similaires, mais également dans quelle mesure ces résultats concordent.

228

229 **Normes.** Une norme correspond à la distribution des scores obtenus par un échantillon
230 de personnes, représentatif d'une population définie, à un instrument qui a été administré dans
231 des conditions standardisées. Tout d'abord, le manuel doit faire état de la population
232 d'étalonnage pour que l'utilisateur puisse savoir si celle-ci est représentative de la situation de
233 son patient. Plusieurs informations sont nécessaires (Ivanova et Hallowell, 2013), ce qui inclut
234 les caractéristiques des enfants formant l'échantillon (p. ex. l'âge et le niveau de scolarité des
235 enfants, la répartition géographique/l'origine, la répartition des statuts socioéconomiques des
236 parents, la proportion de filles et garçons, le nombre d'enfants présentant des difficultés
237 intégrées dans l'échantillon des enfants sans difficulté, etc.). En bref, l'échantillon doit être
238 décrit en précision. La taille de l'échantillon est également une variable importante à

239 considérer. Selon le consensus généralement établi et rapporté dans Gaul Bouchard et al.
240 (2009) et dans Leclercq et Veys (2014), la loi de la limite inférieure exige un minimum de
241 100 personnes dans chaque sous-groupe. Ajoutons qu'il est important que le manuel du test
242 précise le moment de l'étalonnage et les qualifications de l'évaluateur.

243 Enfin, les tests doivent faire état des mesures de tendance centrale (Ivanova et
244 Hallowell, 2013), c'est-à-dire de la moyenne et de l'écart-type (ou des rangs centiles de
245 performances) de l'échantillon d'étalonnage, afin d'avoir un repère quantitatif clair auquel
246 comparer les performances des enfants, et ainsi, être en mesure de les situer par rapport à la
247 moyenne, ou encore, de mettre en évidence leur faiblesse ou leur déficit. Le DSM-5 préconise
248 l'utilisation de tests formels ciblés (c.-à-d. standardisés et normés) avec un seuil de
249 performance correspondant à 1,5 écarts-types sous la moyenne, ou encore, au 7^e percentile
250 pour conclure à la présence d'un trouble des apprentissages en mathématiques. Le DSM-5
251 précise également qu'un seuil plus indulgent (p. ex. -1 écart-type) peut être utilisé pour
252 identifier la présence de difficultés en mathématiques. Abondant dans une direction similaire,
253 Green et Gallagher (2014) rapportent dans leur synthèse de la littérature que la recherche
254 scientifique considère que des scores à des tests évaluant les habiletés mathématiques se
255 situant sous le 10^e percentile indiqueraient la présence d'un trouble des apprentissages en
256 mathématiques (*Mathematic Learning Disabilities*), alors que des scores se situant sous le 35^e
257 percentile indiqueraient simplement la présence de difficultés mathématiques (*Mathematic
258 Difficulties*). Ce seuil plus large a notamment sa place dans une démarche de dépistage des
259 enfants à risque. Dans le cas de la prévention, il est effectivement préférable et plus prudent
260 d'obtenir plus de faux positifs que de faux négatifs. Par ailleurs, un score brut d'un test normé
261 n'est qu'une mesure approximative du score véritable de l'individu. Afin de minimiser
262 l'impact de cette estimation, McCauley et Swisher (1984) préconisent de fixer un intervalle de
263 confiance (IC) à 95%. Il s'agit d'un intervalle de valeurs (dépendant de l'écart-type de la

264 distribution des scores et du degré de fidélité des tests) qui détermine 95% de chance de
265 contenir la vraie valeur du paramètre estimé. Autrement dit, l'intervalle de confiance
266 représente la fourchette de valeurs à l'intérieur de laquelle nous sommes certains à 95% de
267 trouver la vraie valeur recherchée.

268

269 **Outils d'évaluation des capacités mathématiques**

270 Lafay, St-Pierre et Macoir (2014) ont réalisé une recension des outils disponibles en
271 français pour l'évaluation mathématique et ont conclu que les professionnels ont quelques
272 outils à disposition pour évaluer les habiletés mathématiques des enfants. Ces auteurs mettent
273 toutefois en évidence des limites, telles que le manque de standardisation ou de normes pour
274 plusieurs outils, ou encore, le fait que certains outils ne s'appuient pas sur les modèles
275 théoriques actuels du traitement numérique et ne permettent donc pas de documenter les
276 processus déficitaires nécessaire pour diagnostiquer un trouble des apprentissages en
277 mathématiques. Jusqu'à présent, aucune étude ne s'est intéressée à évaluer les qualités
278 psychométriques des tests disponibles en français pour évaluer les capacités mathématiques
279 des enfants francophones.

280

281 **Objectifs**

282 L'objectif général du présent article est d'aider l'orthophoniste à faire un choix éclairé
283 dans la sélection des outils d'évaluation mathématique dont il a besoin en plus de lui
284 permettre d'être en mesure de mieux comprendre ou d'interpréter les résultats des évaluations
285 complétées par d'autres professionnels. Pour cela, les objectifs spécifiques sont : 1) mettre à
286 jour la recension des outils disponibles en français pour l'évaluation mathématique établie par
287 Lafay et al. (2014) et 2) faire une analyse des qualités psychométriques des outils standardisés
288 faisant partie de la recension.

289

290

Recension des outils

291 Méthodologie

292 La recension a porté sur les outils permettant l'évaluation des capacités mathématiques
293 auprès de la population pédiatrique francophone. Celle-ci a d'abord été effectuée à partir des
294 résultats de la recension de Lafay et al. (2014). Une mise à jour a ensuite été effectuée en
295 utilisant plusieurs moyens. Une première recherche a été effectuée dans les bases de données
296 PubMed et PsycInfo à l'aide des mots-clés « évaluation » et « mathématiques ». Une
297 recherche identique a également été réalisée sur le site de la revue orthophonique *Glossa*,
298 puisque celle-ci n'est pas référencée dans les bases de données mentionnées précédemment.
299 Toutefois, la plupart des outils d'évaluation ne sont pas référencés dans les bases de données
300 scientifiques. De ce fait, les catalogues des grandes maisons d'édition de tests (c.-à-d. Édition
301 du centre de psychologie appliquée et Pearson) et des maisons d'édition spécialisées dans le
302 matériel orthophonique (c.-à-d. Orthoédition, HappyNeuron, Orthopratic et Cogilud) ont été
303 consultés. Les catalogues ont été parcourus page à page dans les rubriques concernant
304 l'évaluation et les mathématiques. Finalement, des chercheurs dans le domaine de la
305 psychologie ou de l'éducation en mathématiques, ainsi que des cliniciens (orthophonistes,
306 neuropsychologues et orthopédagogues), ont été consultés dans le but d'identifier d'autres
307 outils utilisés. La recherche a été menée afin de repérer les outils permettant le dépistage et
308 l'évaluation des difficultés mathématiques auprès la population pédiatrique francophone
309 édités entre l'année 1990 et février 2017.

310

311 Résultats

312 La recension de Lafay et al. (2014) avait mené à l'identification de 25 outils : trois
313 échelles d'intelligence, six outils d'évaluation du rendement scolaire et 15 outils spécialisés

314 dans l'évaluation mathématique cognitive. La présente recherche a mené à l'identification de
315 six outils supplémentaires, pour un total de 31 outils. En effet, les quatre outils suivants ont
316 été repérés dans les catalogues de maisons d'édition : Evaluation Des fonctions cognitives et
317 des Apprentissages (EDA; Billard et Touzin, 2012), Épreuve verbale d'aptitudes cognitives
318 (EVAC; Flessas et Lussier, 2003), Tedi-math Grands (Noël et Grégoire, 2015) et Examath 8-
319 15 (Lafay et Helloin, 2016). De plus, deux articles présentant les données de normalisation en
320 franco-qubécois d'outils existants, soit le Zareki-R (Dellatolas et Von Aster, 2006; Lafay, St-
321 Pierre et Macoir, 2016) et le Tempo Test Rekenen (TTR; De Vos, 1992; Lafay, St-Pierre et
322 Macoir, 2015), ont été repérés.

323 Parmi les 31 outils répertoriés, neuf outils ont été retirés à la suite d'application de
324 critères d'exclusion. D'abord, trois échelles d'intelligence ont été retirées (Kaufman
325 Assessment Battery for Children , KABC-II ; Wechsler Intelligence Scale for Children,
326 WISC-IV ; Nouvelle échelle métrique de l'intelligence, NEMI-2) car le sous-test de
327 mathématiques ne peut donner sens sans le contexte complet de l'échelle d'intelligence et car
328 le public lecteur visé est principalement l'orthophoniste qui ne fait pas passer d'échelle
329 d'intelligence. De plus, deux outils ont été retirés parce que le manuel n'était pas disponible
330 pour consultation (BATELEM ; Tests d'acquisitions scolaires mathématiques, TAS), deux
331 autres outils car il s'agissait d'épreuves totalement descriptives (Épreuve de décision logique
332 et Difficultés en mathématiques, évaluation et rééducation). Enfin, deux outils ont été retirés
333 parce que le manuel était rédigé dans une autre langue que le français (Keymath-3 et
334 Kortrijkse Rekestest Revisie, KRT). Cela résulte en un total de 22 outils ayant été inclus dans
335 la présente étude. Parmi ceux-ci, 14 sont des outils évaluant uniquement les habiletés
336 mathématiques, alors que les huit autres sont des batteries de langage ou de rendement
337 scolaire comportant un ou quelques subtests d'évaluation mathématique. Le tableau 1
338 présente les caractéristiques générales des tests : le titre, le(s) auteur(s), la date de publication,

339 les informations sur la modalité de présentation (informatisée ou papier), ainsi que les
340 caractéristiques de la population et du moment d'étalonnage. Les batteries évaluant
341 uniquement les habiletés mathématiques permettent d'évaluer les enfants âgés de 4 ans 0 mois
342 à 17 ans 11 mois. Huit sont normées pour la population française, trois pour la population
343 belge francophone, trois pour la population franco-québécoise et un pour la population suisse
344 francophone.

345

346 **[Insérer ici le Tableau 1]**

347

348 Les outils mesurent les habiletés mathématiques (c.-à-d. dénombrement, numération,
349 transcodage, calcul, vocabulaire mathématique, résolution de problèmes, raisonnement) et les
350 habiletés cognitives de traitement du nombre. Les domaines mathématiques couverts varient
351 toutefois d'un outil à l'autre. Nous renvoyons le lecteur au tableau 2 pour le détail des
352 domaines mathématiques couverts par chaque outil.

353

354 **[Insérer ici le Tableau 2]**

355

356 **Analyse des qualités psychométriques des outils**

357 **Méthodologie**

358 Un total de vingt-deux outils a été soumis à l'analyse des qualités psychométriques.
359 Une grille d'analyse a été construite pour les besoins de l'étude à partir d'une synthèse de
360 plusieurs références traitant du sujet des qualités psychométriques d'outils d'évaluation (c.-à-
361 d. Gaul Bouchard et al., 2009; Ivanova et Hallowell, 2013; Leclercq et Veys, 2014). Ces
362 références ont été choisies pour deux raisons principales : elles portaient sur un domaine
363 proche de celui de la présente étude, à savoir un domaine de l'orthophonie (c.-à-d. le langage

364 oral), et elles apportaient une analyse rigoureuse d'autres outils d'évaluation. Par exemple,
365 Gaul Bouchard et al. (2009) ont utilisé une grille composée de 16 critères tirés des
366 recommandations de McCauley et Swisher (1984). Leclercq et Veys (2014) ont quant à eux
367 employé une grille composée de 13 critères, ceux-ci également tirés des recommandations de
368 McCauley et Swisher (1984). Si certains critères sont communs dans les deux grilles, certains
369 ne sont présents que dans l'une ou l'autre. Finalement, Ivanova et Hallowell (2013) ont décrit
370 certains autres critères supplémentaires, tels que la nécessité d'évaluer la validité de surface,
371 la structure du test par une analyse factorielle, etc. La mise en commun de ces travaux a ainsi
372 mené à l'élaboration de la présente grille. Au total, 21 critères ont été établis comme étant des
373 caractéristiques de base devant être considérées par le clinicien avant d'utiliser un test dans le
374 but de poser un diagnostic ou d'émettre une décision clinique à propos d'une performance
375 d'un enfant à un test. Ces critères sont présentés et expliqués dans le tableau 3.

376 Un score de qualité est attribué à chaque test : [nombre de critères validés / 21 critères
377 au total * 100]. Un score de validation est calculé pour chaque critère : [nombre de tests ayant
378 validé le critère / 22 tests au total * 100].

379 L'un des critères présentés dans le tableau 3 est la validité de contenu. Celle-ci a été
380 établie à partir de la lecture d'ouvrages de référence en cognition mathématique
381 qui font état des lieux des modèles théoriques actuels du développement mathématique et du
382 trouble des apprentissages en mathématiques chez l'enfant (Butterworth, 1999, 2005;
383 Cappelletti et Fias, 2016; Dehaene, 2010; Habib, Noël, George-Poracchia et Brun, 2011;
384 Habib, 2014; Kadosh et Dowker, 2015), afin d'identifier si l'outil d'évaluation mathématique
385 ou les subtests de batteries plus générales s'appuyaient sur les modèles du trouble des
386 apprentissages en mathématiques définis dans ces ouvrages. Dans le cas d'une batterie de
387 langage comportant quelques subtests mathématiques, la validité de contenu a été établie, non
388 pas pour la batterie au complet, mais pour les subtests en question. En particulier, nous avons

389 accordé 1 point à un outil faisant référence à un ouvrage précis et à un modèle actuel basé sur
390 les données probantes, 0,5 point à un outil indiquant une référence mais n'expliquant pas le
391 modèle précis et manquant ainsi de précision mais basée sur des données probantes et 0 point
392 à un outil s'appuyant sur un modèle théorique non reconnu par la littérature actuelle ou à un
393 outil ne précisant aucune référence.

394

395

[Insérer ici le Tableau 3]

396

397 Un test a été considéré comme satisfaisant un critère si le manuel présentait, dans son
398 entier, suffisamment d'informations en lien avec le critère en question pour en permettre
399 l'évaluation. Dans ce cas, 1 point a été attribué. Au contraire, si aucune information n'était
400 donnée, aucun point n'a été attribué. Dans certains cas, nous avons décidé d'accorder
401 seulement 0,5. Par exemple, si le manuel stipulait qu'une des validités considérées dans
402 l'analyse psychométrique avait été vérifiée mais qu'aucune donnée chiffrée ne permettait de
403 réellement approuver la présence du critère, 0,5 point était accordé en guise de confiance aux
404 auteurs. Ce même score (0,5 point) était également attribué si le manuel donnait les
405 informations relatives à un critère, mais les données statistiques révélaient des résultats non
406 significatifs ou faibles (p. ex. une corrélation faible). Enfin, dans trois cas, nous avons décidé
407 d'attribuer 0,75 point, car les manuels indiquaient que les tests remplissaient presque
408 totalement le critère. Par exemple, un score de 0,75 point a été attribué à l'outil Examath 8-15
409 pour le critère « Taille de l'échantillon », car certains groupes dépassaient le seuil de 100
410 enfants alors que d'autres en était proche (p. ex. 87).

411

412

413

Deux types de totaux ont ainsi été calculés. Premièrement, le nombre de critères
remplis pour chaque test a été additionné, chaque critère possédant une importance relative
équivalente dans la présente grille constituée (c.-à-d. 1 critère = 1 point). Toutefois, dans les

414 faits, certains critères nous semblent avoir une plus grande importance (p. ex., la
415 sensibilité/spécificité versus la validité d'apparence). Une note totale sur 21 a été attribuée et
416 le pourcentage correspondant calculé. Deuxièmement, le nombre de tests remplissant chaque
417 critère a été additionné, chaque critère possédant une importance relative équivalente (c.-à-d.
418 1 test = 1 point). Une note totale sur 22 a été attribuée et le pourcentage correspondant
419 calculé.

420 L'analyse psychométrique a été effectuée par deux juges (auteurs de l'article) qui
421 sont toutes deux orthophonistes et impliquées dans la recherche et/ou la pratique
422 orthophonique basée sur les données probantes. Chacune a suivi une formation de base sur les
423 qualités psychométriques des outils d'évaluation dans son cursus de formation continue
424 professionnelle. Tout d'abord, les deux juges ont évalué chaque outil, de manière séparée et à
425 l'aveugle, à partir de la consultation des manuels, de la consultation des sites commerciaux et
426 des échanges avec les auteurs (quand ceux-ci ont bien accepté de répondre aux interrogations
427 suscitées par la lecture des manuels). Ensuite, le premier juge a vérifié l'adéquation des points
428 attribués par elle-même et son co-juge (et inversement). L'analyse a d'abord montré une
429 adéquation globale de 87% entre les deux juges. Ainsi, les deux juges ont discuté et parcouru
430 à nouveau les manuels ensemble pour parvenir à un consensus complet. C'est d'ailleurs à ce
431 moment qu'ont été définies les cotations intermédiaires de 0,5 ou 0,75 précédemment
432 explicitées.

433

434 **Résultats**

435 Les tableaux 4a et 4b présentent une synthèse des caractéristiques psychométriques
436 des 22 outils analysés : la note de 1 point, 0,75 point, 0,50 point ou 0 point est indiquée dans
437 chaque case.

463 permettre d'être en mesure de mieux comprendre ou d'interpréter les résultats d'évaluations
464 complétées par d'autres professionnels. Pour cela, les objectifs spécifiques étaient de : 1)
465 mettre à jour la recension des outils disponibles en français pour l'évaluation mathématique
466 établie par Lafay et al. (2014) et 2) faire une analyse des qualités psychométriques des outils
467 standardisés faisant partie de la recension.

468

469 **Premier objectif : recension des outils**

470 Relativement au premier objectif, la présente étude a permis de mettre en évidence
471 l'existence de 22 outils disponibles en français pour l'évaluation mathématique, dont 14 outils
472 évaluant spécifiquement les habiletés mathématiques et huit batteries de langage ou de
473 rendement scolaire comportant un ou quelques subtests d'évaluation mathématique. Les outils
474 évaluant spécifiquement les habiletés mathématiques permettent d'évaluer les enfants âgés de
475 4 ans 0 mois à 17 ans et 11 mois. Huit sont normées pour la population française, trois pour la
476 population belge francophone, trois pour la population franco-québécoise et un pour la
477 population suisse francophone. Les domaines mathématiques couverts sont les habiletés
478 mathématiques (dénombrement, numération, transcodage, calcul, vocabulaire mathématique,
479 résolution de problèmes, raisonnement) et les habiletés cognitives de traitement du nombre.
480 Ils varient toutefois d'un outil à un autre.

481

482 **Deuxième objectif : analyse des qualités psychométriques des outils**

483 Concernant les qualités psychométriques des outils d'évaluation mathématique
484 recensés, la présente étude a permis de mettre en évidence que certains critères sont très bien
485 considérés alors que d'autres ne le sont pas ou quasiment. De plus, l'analyse a montré que les
486 outils n'ont pas tous un score de qualité psychométrique global équivalent.

487 Tout d'abord, l'analyse montre qu'aucun critère n'est pris en compte par l'ensemble
488 des tests. De manière générale, les outils que les professionnels ont à leur disposition
489 respectent plutôt bien les critères suivants : la qualification de l'évaluateur, la standardisation
490 de l'outil (consigne de passation et cotation), la précision de l'objectif de l'épreuve, la
491 description de l'échantillon d'étalonnage et la présence de mesures de tendance centrale. Ces
492 critères sont les plus simples à considérer et à mettre en place d'après Gaul Bouchard et al.
493 (2009). En revanche, neuf des 21 critères ont été presque totalement négligés dans les tests
494 disponibles en français, soit la validité de surface (ou d'apparence), la validité de critère
495 concomitante, la conception factorielle de l'outil, la sensibilité et la spécificité de l'outil, la
496 fidélité de type stabilité (temporelle, versions parallèles, inter-juges) et la cohérence interne
497 (bissection et consistance inter-items). En particulier, les outils actuels ont généralement une
498 bonne spécificité mais la sensibilité est limitée. Cela peut s'expliquer car aucune mesure n'est
499 réalisée avec un groupe conséquent d'enfants en difficulté. Or, des scores-seuils devraient être
500 calculés avec la distribution d'une population de personnes saines et de personnes en
501 difficulté. L'outil Examath 8-15 est le seul outil, dans cette recension, à fournir des données
502 concernant la sensibilité et la spécificité de la batterie. En effet, le manuel présente l'étude des
503 réussites de 126 enfants dont 63 présentaient des difficultés mathématiques et 63 faisaient
504 partie de la population des personnes saines. Leclercq et Veys (2014) déploraient aussi
505 l'absence d'un « critère diagnostique crucial, le pouvoir discriminant des outils » lors de
506 l'analyse des outils d'évaluation du langage pour la population francophone. Le manque
507 d'informations ou de moyens des concepteurs de tests sont des causes possibles à ces
508 absences.

509 De plus, les qualités psychométriques ne sont pas identiques pour tous les outils. Le
510 constat est identique à celui fait par Leclercq et Veys (2014) lors de l'analyse d'outils
511 d'évaluation du langage pour la population francophone : les outils diagnostiques à la

512 disposition des professionnels ne rencontrent pas l'ensemble des critères psychométriques
513 recommandés pour une pratique de qualité. En effet, la présente analyse montre qu'aucun test
514 n'obtient un score de qualité de 100%, ou encore, au-dessus de 75%. Si certains auteurs ont
515 fait de gros efforts quant aux qualités psychométriques de leur outil (p. ex. les auteurs de
516 l'Examath 8-15, du Tedi-math Grands, du Wechsler Individual Achievement Test (WIAT-II),
517 de l'Exalang 8-11 et de l'Exalang 11-15), peu d'outils obtiennent un score de qualité
518 supérieur à 50%. Ce résultat mène à une recommandation de privilégier les trois outils
519 évaluant spécifiquement les habiletés mathématiques qui obtiennent un score de qualité
520 supérieur à 50% (c.-à-d. l'Examath 8-15, le Tedi-math Grands et le WIAT-II).

521 La présente étude met ainsi en lumière l'existence d'un écart considérable entre les
522 qualités psychométriques des outils d'évaluation présentées et celles souhaitées. Elle permet,
523 en cela, d'aider le clinicien à faire un choix éclairé dans la sélection des outils dont il a besoin,
524 dans la limite des choix existants. Plusieurs critères sont primordiaux à considérer : la
525 standardisation, la validité, la fidélité et les caractéristiques de normalisation. Selon Gaul
526 Bouchard et al. (2009), « *puisque ces tests sont utilisés par des professionnels de différents*
527 *domaines, ces critères sont importants car ils assurent d'obtenir des informations plus*
528 *objectives. Conséquemment, les décisions cliniques qui en découlent risquent moins d'être*
529 *influencées par la manière dont les praticiens conceptualisent et interprètent les construits*
530 *évalués* ». La validité de contenu est particulièrement importante dans la mesure où il ne fait
531 pas sens, dans une démarche de pratique basée sur les données probantes, d'utiliser un test
532 employant des tâches reposant sur un modèle théorique invalide. De même, le pouvoir
533 discriminant (autrement dit la sensibilité d'un test) paraît des plus indispensables, puisque
534 c'est la qualité qui permet d'attester qu'un outil d'évaluation permet de repérer l'ensemble des
535 enfants présentant un trouble des apprentissages en mathématiques.

536

537 Limites

538 L'étude s'est confrontée à trois limites principales. La première concernait la difficulté
539 à retracer les outils disponibles, car très peu ont fait l'objet de publications scientifiques et ils
540 n'étaient donc pas recensés dans les bases de données scientifiques. La deuxième concernait
541 la difficulté à retrouver les informations dans les manuels, car ceux-ci ne possédaient pas tous
542 la même structure et n'utilisaient pas nécessairement le même vocabulaire. De plus, si le
543 constat est fait que tous les outils ne remplissent pas toutes les qualités psychométriques
544 évaluées dans la présente étude, il convient de nuancer le propos. En effet, la qualité est en
545 fonction des critères retenus dans cet article. Il ne s'agit pas d'une valeur absolue, mais bien
546 d'un résultat en fonction de la grille élaborée pour les besoins de la présente étude. Celle-ci
547 donne une indication relative. Enfin, une autre limite de l'étude concernait l'importance
548 relative équivalente des critères évalués (tel que discuté précédemment). En effet, deux tests
549 pourraient obtenir un résultat équivalent en pourcentage, mais détenir des caractéristiques
550 psychométriques bien différentes (dont certaines pourraient être plus importantes que d'autres
551 lorsque vient le temps de choisir un test). Il semble alors indispensable de ne pas uniquement
552 s'en tenir au résultat final pour caractériser les qualités psychométriques d'un outil et pour
553 faire un choix d'outil, mais bien d'analyser l'ensemble des critères.

554

555 Conclusion

556 Jusqu'à présent, aucune étude ne s'était intéressée à évaluer les qualités
557 psychométriques des tests disponibles en français pour l'évaluation des capacités
558 mathématiques des enfants francophones. La présente étude est donc tout à fait originale et
559 pertinente dans le contexte de la pratique orthophonique basée sur les données probantes.
560 Vingt-deux outils ont été recensés et leurs propriétés psychométriques ont été analysées.
561 L'étude a mis en évidence le fait que certains critères sont très bien considérés (p. ex. la

562 standardisation) alors que d'autres ne le sont pas ou quasiment (p. ex. le pouvoir
563 discriminant). De plus, tous les outils n'ont pas un score de qualité psychométrique global
564 équivalent. Parmi les 22 recensés, seulement trois outils évaluant spécifiquement les habiletés
565 mathématiques obtiennent un score de qualité supérieur à 50% (c.-à-d. l'Examath 8-15, le
566 Tedi-math Grands et le WIAT-II). Il faut toutefois noter que tout critère n'a pas la même
567 importance. Il semble alors indispensable de ne pas s'en tenir uniquement au score global
568 pour caractériser les qualités psychométriques d'un outil et pour faire un choix d'outil, mais
569 bien d'analyser l'ensemble des critères.

570 Quelques recommandations générales peuvent être développées. À l'avenir, il est
571 indispensable que les futurs concepteurs d'outils d'évaluation mathématique fassent l'effort
572 de développer des outils standards, d'investiguer la validité et la fidélité des outils et de
573 donner un maximum de précision quant à l'échantillon d'étalonnage et les normes dans les
574 manuels d'utilisation, pour une plus grande transparence. De plus, tout comme il existe un
575 canevas général tacitement accepté et utilisé pour la rédaction d'articles scientifiques, un
576 canevas général concernant la rédaction de manuels de tests devrait être développé et utilisé
577 par les concepteurs et les maisons d'édition. Chaque manuel devrait ainsi détailler les aspects
578 liés à 1) la standardisation, 2) la validité, 3) la fidélité et 4) la normalisation de l'outil.

579 De même, il est indispensable que les cliniciens considèrent l'ensemble de ces critères
580 pour juger des outils valides et pertinents à utiliser (Betz et al., 2013 ; Gaul Bouchard et al.,
581 2009; Leclercq et Veys, 2014; McCauley, 1989). Néanmoins, l'obstacle principal à la mise en
582 place d'une pratique basée sur les données probantes, d'après les informations recueillies
583 auprès d'orthophonistes provenant de différents pays, reste le manque de temps (Durieux et
584 al., 2016; O'Connort et Pettigrew, 2009; Zipoli et Kennedy, 2005). Aider les cliniciens à
585 analyser les tests est donc essentiel pour que ces derniers soient conscientisés à la
586 répercussion de l'absence de certaines qualités psychométriques sur leur pratique clinique.

587 **Références**

- 588 American Psychiatric Association. (2013). *Diagnostic and statistical manual of mental disorders DSM-*
589 *5* (5e éd.). Arlington, VA : American Psychiatric Publishing. American Psychiatric Association.
590 (2016). *DSM-5 : manuel diagnostique et statistique des troubles mentaux* (5e éd.). Issy-les-
591 Moulinaux, France : Elsevier Masson.
- 592 Antell, S. E. et Keating, D. P. (1983). Perception of numerical invariance in neonates. *Child*
593 *Development*, 54, 695–701. doi:10.2307/1130057
- 594 Barbaresi, W. J., Katusic, S. K., Colligan, R. C., Weaver, A. L. et Jacobsen, S. J. (2005). Math
595 learning disorder: Incidence in a population-based birth cohort, 1976-82, Rochester,
596 Minn. *Ambulatory Pediatrics*, 5, 281–289. doi:10.1367/A04-209R.1
- 597 Baudonck, M., Debusschere, A., Dewulf, B., Samyn, F. et Vercaemst, V. (2006). *KRT-R*
598 *2006. Kortrijkse Rekentest Revisie ou test de calcul de Courtrai révisé*. Courtrai,
599 Belgique : Revalidatiecentrum Overleie.
- 600 Betz, S. K., Eickhoff, J. R. et Sullivan, S. F. (2013). Factors influencing the selection of
601 standardized tests for the diagnosis of specific language impairment. *Language, Speech,*
602 *and Hearing Services in Schools*, 44, 133–46. doi:10.1044/0161-1461(2012/12-0093)
- 603 Billard, C. et Touzin, M. (2012). *Evaluation Des fonctions cognitives et des Apprentissages*
604 *de 4 à 11 ans*. Isbergues, France : Ortho Édition.
- 605 Brannon, E. M. (2005). What animals know about numbers. Dans J. I. D. Campbell (dir.),
606 *Handbook of mathematical cognition* (p. 85–107). New-York, NY : Psychology Press.
- 607 Butterworth, B. (1999). *The mathematical brain*. London, United Kingdom : MacMillan.

- 608 Butterworth, B. (2005). The developmental dyscalculia. Dans J. I. D. Campbell (dir.),
609 *Handbook of mathematical cognition*, (p. 455–467). New-York, NY : Psychology Press.
- 610 Butterworth, B., Reeve, R., Reynolds, F. et Lloyd, D. (2008). Numerical thought with and
611 without words: Evidence from indigenous Australian children. *Proceedings of the*
612 *National Academy of Sciences of the United States of America*, 105, 13 179–13 184.
613 doi:10.1073/pnas.0806045105
- 614 Cappelletti, M. et Fias, W. (2016). *Progress in brain research: Vol. 227. The mathematical*
615 *brain across the lifespan*. Amsterdam, Pays-Bas : Elsevier.
- 616 Carey, S. (2001). Cognitive foundations of arithmetic: Evolution and ontogenesis. *Mind &*
617 *Language*, 16, 37–55. doi:10.1111/1468-0017.00155
- 618 Carey, S. (2004). Bootstrapping & the origin of concepts. *Doedalus*, 131(1), 59–68.
619 doi:10.1162/001152604772746701
- 620 Chevrie-Muller, C. et Plaza, M. (2001). *N-EEL. Nouvelles Épreuves pour l'Examen du*
621 *Langage*. Paris, France : Édition du Centre de Psychologie Appliquée.
- 622 Cagnet, G. (2006). *NEMI-2. Nouvelle échelle métrique de l'intelligence-2*. Paris, France :
623 Édition du Centre de Psychologie Appliquée.
- 624 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale,
625 NJ : L. Erlbaum Associates.
- 626 Connoly, A. J. (2008). *KeyMathTM 3 diagnostic assessment: Canadian edition*. San Antonio,
627 TX : Pearson.

- 628 de Barbot, F., Duquesne, F., Marchand, M.H., Mazeau, M., Meljac, C., Truscelli, D.,
629 Vergnaud, G. (1995). ECPN. Epreuves Conceptuelles de résolution des Problèmes
630 Numériques. CIMETE (non édité).
- 631 Dehaene, S. (2010). *La bosse des maths. 15 ans après*. Paris, France : Odile Jacob.
- 632 De Vos, T. (1992). *Tempo Test Rekenen*. Berkhout, Pays-bas : Nijmegen.
- 633 Donlan, C., Cowan, R., Newton, E. J. et Lloyd, D. (2007). The role of language in
634 mathematical development: Evidence from children with specific language impairments.
635 *Cognition*, 103, 23–33. doi:10.1016/j.cognition.2006.02.007
- 636 Durieux, N., Pasleau, F., Piazza, A., Donneau, A.-F., Vandenput, S. et Maillart, C. (2016).
637 Information behaviour of French-speaking speech-language therapists in Belgium:
638 Results of a questionnaire survey. *Health Information and Libraries Journal*, 33, 61–76.
639 doi:10.1111/hir.12118
- 640 Durkin, K., Mok, P. L. H. et Conti-Ramsden, G. (2013). Severity of specific language
641 impairment predicts delayed development in number skills. *Frontiers in Psychology*,
642 4(581), 1–10. doi:10.3389/fpsyg.2013.00581
- 643 Flessas, J. et Lussier, F. (2003). *EVAC, Épreuve verbale d'aptitudes cognitives*. Paris, France :
644 Édition du Centre de Psychologie Appliquée.
- 645 Fleiss, J. L. (1981) *Statistical methods for rates and proportions*. 2nd ed. (New York: John
646 Wiley)

- 647 Friberg, J. C. (2010). Considerations for test selection: How do validity and reliability impact
648 diagnostic decisions? *Child Language Teaching and Therapy*, 26(1), 77–92.
649 doi:10.1177/0265659009349972
- 650 Frank, M. C., Everett, D. L., Fedorenko, E. et Gibson, E. (2008). Number as a cognitive
651 technology: Evidence from Pirahã language and cognition. *Cognition*, 108, 819–824.
652 doi:10.1016/j.cognition.2008.04.007
- 653 Gaillard, F. (2000). *Numerical : test neurocognitif pour l'apprentissage du nombre et du*
654 *calcul*. Lausanne, Suisse : Institut de psychologie Université de Lausanne.
- 655 Gaul Bouchard, M.-E., Fitzpatrick, E. M. et Olds, J. (2009). Analyse psychométrique d'outils
656 d'évaluation utilisés auprès des enfants francophones. *Revue canadienne d'orthophonie*
657 *et d'audiologie*, 33, 129–139.
- 658 Gordon, P. (2004). Numerical cognition without words: Evidence from Amazonia. *Science*,
659 306, 496–499. doi:10.1126/science.1094492
- 660 Green, K. B. et Gallagher, P. A. (2014). Mathematics for young children: A review of the
661 literature with implications for children with disabilities. *Başkent University Journal of*
662 *Education*, 1(1), 81–92.
- 663 Gross-Tsur, V., Manor, O. et Shalev, R. S. (1996). Developmental dyscalculia: Prevalence
664 and demographic features. *Developmental Medicine & Child Neurology*, 38, 25–33.
665 doi:10.1111/j.1469-8749.1996.tb15029.x
- 666 Habib, M. (2014). *La Constellation des dys*. Paris, France : De Boeck-Solal.

- 667 Habib, M., Noël, M.-P., George-Poracchia, F. et Brun, V. (2011). *Calcul et dyscalculie : des*
668 *modèles à la rééducation*. Paris, France : Masson. Heremans, M. (2011). *MathEval*
669 *Dépistage de la dyscalculie*. Repéré à <https://sites.google.com/site/testmatheval/>
- 670 Ivanova, M. V. et Hallowell, B. (2013). A tutorial on aphasia test development in any
671 language: Key substantive and psychometric considerations. *Aphasiology*, 27, 891–920.
672 doi:10.1080/02687038.2013.805728
- 673 Kadosh, R. C. et Dowker, A. (2015). *The Oxford handbook of numerical cognition*. Oxford,
674 United Kingdom: Oxford library of psychology.
- 675 Kaufmann, A. S. et Kaufmann, N. L. (2008). *KABC-II. Batterie pour l'examen psychologique*
676 *de l'enfant - 2^{ème} édition*. Paris, France : Édition du Centre de Psychologie Appliquée.
- 677 Lafay, A. et Helloin, M.-C. (2016). *Examath 8-15 : batterie informatisée d'examen des*
678 *habiletés mathématiques*. Grenade, France : HappyNeuron.
- 679 Lafay, A., St-Pierre, M.-C. et Macoir, J. (2014). L'évaluation des habiletés mathématiques de
680 l'enfant : inventaire critique des outils disponibles. *Glossa*, 116, 33–58.
- 681 Lafay, A., St-Pierre, M.-C. et Macoir, J. (2015). Validation franco-qubécoise du Tempo Test
682 Rekenen pour l'évaluation des habiletés mathématiques auprès d'enfants de 8-9 ans.
683 *Glossa*, 118, 27–39.
- 684 Lafay, A., St-Pierre, M.-C. et Macoir, J. (2016). Performances moyennes des enfants franco-
685 québécois de 8-9 ans au test mathématique Zareki-R. *Glossa*, 119, 41–54.

- 686 Leclercq, L. et Veys, E. (2014). Réflexions sur le choix de tests standardisés lors du
687 diagnostic de dysphasie. *Approche Neuropsychologique des Apprentissages chez*
688 *l'Enfant*, 26. 374–382
- 689 Legeay, M. P., Morel, L. et Voye, M. (2009). *Mallette ERLA (Exploration du Raisonnement*
690 *et du Langage Associé)*. Trucy sur Yonne, France : Cogilud.
- 691 McCauley, R. J. (1989). Measurement as a dangerous activity. *Revue d'orthophonie et*
692 *d'audiologie*, 13, 29–32.
- 693 McCauley, R. J. et Swisher, L. (1984). Use and misuse of norm-referenced tests in clinical
694 assessment: A hypothetical case. *Journal of Speech and Hearing Disorders*, 49, 338–
695 348. doi:10.1044/jshd.4904.338
- 696 Meljac, C. et Lemmel, G. (1999). *UDN 2 Construction et Utilisation du nombre*. Paris, France
697 : Édition du Centre de Psychologie Appliquée.
- 698 Ménessier, A. (2003). Les variations stratégiques chez l'enfant dans le calcul d'additions et de
699 soustractions élémentaires. *Glossa*, 83, 20–33.
- 700 Métral, E. (2008). *Mallette B-LM cycle II*. Chavanod, France : Orthopratic.
- 701 Noël, M.-P. et Grégoire, J. (2015). *Tedi-math Grands*. Paris, France : Édition du Centre de
702 Psychologie Appliquée.
- 703 Noël, M.-P. et Rousselle, L. (2011). Developmental changes in the profiles of dyscalculia: An
704 explanation based on a double exact-and-approximate number representation model.
705 *Frontiers in Human Neuroscience*, 5(165), 1–4. doi:10.3389/fnhum.2011.00165
- 706 Nunnally, J. C. (1978). *Psychometric theory* (2nd ed.). New York, NY : McGraw-Hill.

- 707 O'Connor, S. et Pettigrew, C. M. (2009). The barriers perceived to prevent the successful
708 implementation of evidence-based practice by speech and language therapists.
709 *International Journal of Language & Communication Disorders*, 44, 1018–1035.
710 doi:10.1080/13682820802585967
- 711 Plante, E. et Vance, R. (1994). Selection of preschool language test: A data-based approach.
712 *Language, Speech, and Hearing Services in School*, 25, 15–24. doi:10.1044/0161-
713 1461.2501.15
- 714 Riquier, M. (1997). *TAS Révisés – Tests d'acquisitions scolaires mathématiques*. Paris,
715 France : Édition du Centre de Psychologie Appliquée.
- 716 Roux, M.-O. (2014). Surdit  et difficult s d'apprentissage en mathématiques,  tat des lieux et
717 probl matiques actuelles. *Bulletin de Psychologie*, 4, 295–307.
718 doi:10.3917/bupsy.532.0295
- 719 Savigny, M. (2001). *BATELEM-R Batterie d' preuves pour l' cole  l mentaire*. Paris, France
720 :  dition du Centre de Psychologie Appliqu e.
- 721 Simonart, G. (1998a). *ECHAS.  CHelle des Apprentissages Scolaires*. Braine-le-ch teau,
722 Belgique : Eurotests  dition.
- 723 Simonart, G. (1998b). *PEDA1C. Tests p dagogiques de premier cycle primaire*. Braine-le-
724 ch teau, Belgique : Eurotests  dition.
- 725 Starkey, P. et Cooper, R. G. (1980). Perception of numbers by human infants. *Science*, 210,
726 1033–1035. doi:10.1126/science.7434014

- 727 Thibault, M.-P. et Helloin, M.-C. (2006). *Exalang 3-6 : batterie d'examen des fonctions*
728 *langagières chez l'enfant de 3 à 6 ans*. Mont-Saint-Aignan, France : Orthomotus
- 729 Thibault, M.-P., Helloin, M.-C. et Lenfant, M. (2009). *Exalang 11-15 : batterie informatisée*
730 *pour l'examen du langage oral, du langage écrit et des compétences transversales chez*
731 *le collégien*. Mont-Saint-Aignan, France : Orthomotus.
- 732 Thibault, M.-P., Lenfant, M. et Helloin, M.-C. (2012). *Exalang 8-11 : bilan informatisé pour*
733 *l'examen du langage et des compétences transversales chez l'enfant de 8 à 11 ans*.
734 Mont-Saint-Aignan, France : Orthomotus.
- 735 Van Nieuwenhoven, C., Grégoire, J. et Noël, M.-P. (2001). *Tedi-Math. Test diagnostique des*
736 *compétences de base en mathématiques*. Paris, France : Édition du Centre de
737 Psychologie Appliquée.
- 738 Von Aster, M. G. et Shalev, R. S. (2007). Number development and developmental
739 dyscalculia. *Developmental Medicine & Child Neurology*, 49, 868–873.
740 doi:10.1111/j.1469-8749.2007.00868.x
- 741 Von Aster, M. (2006). *Zareki-R : batterie pour l'évaluation du traitement des nombres et du*
742 *calcul chez l'enfant* (adapté par G. Dellatolas). et Paris, France : Édition du Centre de
743 Psychologie Appliquée.
- 744 Wechsler, D. (2005a). *Wechsler Individual Achievement Test - Second edition (WIAT-II)*.
745 London, Angleterre : The Psychological Corporation.
- 746 Wechsler, D. (2005b). *Échelle d'intelligence de Wechsler pour enfants et adolescents -*
747 *quatrième édition (WISC-IV)*. Paris, France : Édition du Centre de Psychologie
748 Appliquée.

749 Wilson, A. J. et Dehaene, S. (2007). Number sense and developmental dyscalculia. Dans D.
750 Coch, G. Dawson et K. W. Fischer (dir.). *Human behavior learning, and the developing*
751 *brain: Atypical development* (p. 212–238). New-York, NY : Guilford Press.

752 Wynn, K. (1992). Addition and subtraction by human infants. *Nature*, 358, 749–750.
753 doi:10.1038/358749a0

754 Zipoli, R. P. et Kennedy, M. (2005). Evidence-based practice among speech-language
755 pathologists: Attitudes, utilization, and barriers. *American Journal of Speech-Language*
756 *Pathology*, 14, 208–220. doi:10.1044/1058-0360(2005/021)

757

Version PREPRINT

Tableau 1

*Relevé des caractéristiques concernant la population d'étalonnage pour les différentes batteries d'évaluation analysées**

Nom du test ou de la batterie	Auteurs	Date	Présentation informatisée	Age (années)	Échantillon : nombre total	Nombre de groupe	Nombre par groupe	Pays	Comparaison avec autres pays	Niveau socioéconomique	Moment de l'étalonnage	
Batteries ou tests mathématiques												
B-LM	Métral	2008	Non	5–8	GSM–CE1	299	14	11 à 31	France	/	/	Janvier-avril
ECPN	Duquesne	2003	Non	4–9	/	132	5	Non indiqué (26 en moyenne)	France	/	/	/
ERLA	Legeay, Morel et Voye	2009	Non	/	/	/	/	/	/	/	/	/
Examath 8–15	Lafay et Helloin	2016	Oui	8–15	CE2–3 ^e collège	443	5	74 à 127	France et Belgique	Québec	Répartition proche des indices INSEE	Mars-mai
MathEval	Heremans	2011	Oui	/	3 ^e maternelle–2 ^e primaire	65	3	65 par groupe environ	Belgique	/	Selon l'auteur, échantillon plutôt favorisé, non représentatif des normes INSEE	Avril à juin
Numerical	Gaillard	2000	Non	7–10	2 ^e –4 ^e primaire	280	2	126 à 154	Suisse	Comparaison entre France, Finlande et Argentine	/	Février
Protocole du calcul élémentaire	Ménissier	2003	Non	7–11	CE1–CM2	406	4	91 à 109	France	/	/	Octobre-novembre
Tedi-math	Van Nieuwenhoven, Grégoire et Noël	2001	Non	/	MSM–CE2 / 2 ^e maternelle–3 ^e primaire	583	8	67 à 76	France et Belgique	Le manuel indique qu'il n'y a pas de différence entre France et Belgique mais données absentes	/	Novembre et mai
Tedi-math Grands	Noël et Grégoire	2015	Les deux	/	CE2–5 ^e collège	254	5	46 à 56	France	/	INSEE	Mars-juin
TTR	Lafay, St-Pierre et Macoir	2015	Non	8–9	3 ^e primaire	77	1	77	Canada francophone (Québec)	Comparaison qualitative avec échantillon du manuel Pays-Bas	Répartition selon Indices de défavorisation du MELS	Février-juin
UDN 2	Meijac et Lemmel	1999	Non	4–11	/	420	8	49 à 57	France	/	/	/
WIAT-II	Wechsler	2005a	Non	6–17 ans 11 mois	1 ^e primaire–5 ^e pecondaire	294 si âge, 304 si classe	12	44 à 56 si classe, 18 à 32 si âge	Canada francophone (Québec)	Le manuel indique une comparaison Franco-	Information sur le niveau d'études des parents	Février à décembre

											Québécois et Franco-Ontariens mais données absentes		
Zareki-R	Dellatolas et Von Aster	2006	Non	6–11 ans et demi	CP–CM2	249	5	43 à 59	France	/	50% ZEP, 28% ne parlent pas français à la maison	Janvier à mars	
Zareki-R	Lafay, St-Pierre et Macoir	2016	Non	8–9	3 ^e Primaire	81	1	81	Canada francophone (Québec)	Comparaison qualitative avec échantillon du manuel France	Répartition selon Indices de défavorisation du MELS	Février-juin	
Batteries générales comportant un ou plusieurs subtests mathématiques													
ECHAS	Simonart	1998a	Non	/	3e–6e primaire	1013	4	201 à 327	Belgique	/	/	Mai-juin	
EDA	Billard et Touzin	2012	Non	4–11 (6–11 pour math)	MSM–CM2 (CP–CM2 pour math)	626	6	94 à 111	France	/	Répartition homogène selon indices INSEE	Septembre-juin	
EVAC	Flessas et Lussier	2003	Non	8–14	CE2–3e collège	886 à 919	7	109 à 154 si âge, 113 à 143 si classe	France	Québec	/	Premier trimestre de l'année	
Exalang 3–6	Thibault et Helloin	2006	Oui	2 ans 8 mois–5 ans 10 mois	MSM–GSM / 2 ^e -3 ^e maternelle	468	6	59 à 96	France et Belgique	/	Répartition proche des indices INSEE	/	
Exalang 8–11	Thibault, Lenfant et Helloin	2012	Oui	8–11	CE2–CM2	461	3	93 à 150	France	/	Répartition proche des indices INSEE	Février-avril	
Exalang 11–15	Thibault, Helloin et Lenfant	2009	Oui	11–15	6 ^e collège–3 ^e collège	322	4	85 à 97	France	/	Répartition proche des indices INSEE	Janvier-mai	
N-EEL	Chevrie-Muller et Plaza	2001	Non	3 ans 7 mois–8 ans 7 mois	PSM–CE2	541	5	108 à 109	France	/	Information sur la catégorie socioprofessionnelle des parents	Septembre-juin	
PEDA 1C	Simonart	1998b	Non	/	1 ^e –2 ^e primaire	232 à 290	3	232 à 290	Belgique	/	/	Mai-juin	

Notes. *Les tests sont classés en ordre alphabétique selon leur titre. **Équivalent des classes entre la France, la Belgique, la Suisse et le Québec : MSM en France = 2^e année de maternelle en Belgique, 1^e année de maternelle en Suisse et prématernelle au Québec; GSM en France = 3^e année de maternelle en Belgique, 2^e année de maternelle en Suisse et maternelle au Québec; CP en France = 1^e année du primaire en Belgique, en Suisse et au Québec; CE1 en France = 2^e année du primaire en Belgique, en Suisse et au Québec; CE2 en France = 3^e année du primaire en Belgique, en Suisse et au Québec; CM1 en France = 4^e année du primaire en Belgique, en Suisse et au Québec; CM2 en France = 5^e année du primaire en Belgique et au Québec, 5^e année de transition en Suisse; 6^e collège en France = 6^e année du primaire en Belgique et au Québec et 6^e année de transition en Suisse; 5^e collège en France : 1^e secondaire en Belgique, en Suisse et au Québec; 4^e collège en France = 2^e secondaire en Belgique, en Suisse et au Québec; 3^e collège en France = 3^e secondaire en Belgique, en Suisse et au Québec; Seconde Lycée en France = 4^e secondaire en Belgique et au Québec et 1^e secondaire degré 2 en Suisse; Première Lycée en France : 5^e secondaire en Belgique et au Québec et 2^e secondaire degré 2 en Suisse. B-LM = Bilan logico-mathématique, cycle 2; CE1 = Cours élémentaire 1; CE2 = Cours élémentaire 2; CM1 = Cours moyen 1; CM2 = Cours moyen 2; CP = Cours préparatoire; ECHAS = ÉCHelle des Apprentissages Scolaires; ECPN = Epreuves Conceptuelles de résolution des Problèmes Numériques; EDA = Évaluation Des fonctions cognitives et Apprentissages; ERLA = Batterie d'Exploration du Raisonnement et du Langage Associé; EVAC = Épreuves Verbale d'Aptitudes Cognitives; Exalang 3–6 = bilan informatisé pour l'examen du langage et des compétences transversales chez l'enfant de 3 à 6 ans; Exalang 8–11 = Bilan informatisé pour l'examen du langage et des compétences transversales chez l'enfant de 8 à 11 ans; Exalang 11–15 = Batterie informatisée pour l'examen du langage oral, du langage écrit et des compétences

transversales chez le collégien; Examath 8–15 = Batterie informatisée d'EXAmen des habiletés MATHématiques; GSM = Grande section de maternelle; INSEE = Institut national de la statistique et des études économiques; MELS = Ministère de l'éducation et de l'enseignement supérieur; MSM = Moyenne section de maternelle; N-EEL = Nouvelles Epreuves d'Examen du Langage; PEDA 1C = Tests pédagogiques de premier cycle primaire; PSM = Petite Section de Maternelle; TTR = Tempo Test Rekenen; UDN 2 = Utilisation du nombre, version 2; WIAT-II = Wechsler Individual Achievement Test, 2ème édition; ZAREKI = Die Neuropsychologische Testbatterie für ZAHlenarbeit und REchnen bei KIndern; ZEP = Zone d'éducation prioritaire

Version PREPRINT

Tableau 2

*Domaines évalués par les différentes batteries d'évaluation analysées**

Nom du test ou de la batterie	Traitement cognitif du nombre	Dénombrement	Numération et transcodage	Calcul	Résolution de problèmes	Langage et raisonnement
Batteries ou tests mathématiques						
B-LM	+/-	+	+	+	+	+
ECPN		+				
ERLA		+	+ (transcodage)		+	+
Examath 8-15	+	+	+	+	+	+
MathEval	+	+	+ (transcodage)	+		
Numerical	+	+	+ (transcodage)	+		+ (vocabulaire math)
Protocole du calcul élémentaire				+		
Tedi-math	+	+	+	+	+	
Tedi-math Grands	+		+	+	+	+
TTR				+		
UDN 2			+ (transcodage)	+		+ (vocabulaire math)
WIAT-II				+	+	+
Zareki-R	+	+	+ (transcodage)	+	+	
Batteries générales comportant un ou plusieurs subtests mathématiques						
ECHAS				+	+	+ (vocabulaire math)
EDA	+	+	+	+	+	+
EVAC						+
Exalang 3-6		+				+ (vocabulaire math)
Exalang 8-11					+	+ (vocabulaire math)
Exalang 11-15						+
N-EEL						+ (vocabulaire math)
PEDA 1C				+		

Notes. *Les tests sont classés en ordre alphabétique selon leur titre. B-LM = Bilan logico-mathématique, cycle 2; ECHAS = Échelle des Apprentissages Scolaires; ECPN = Épreuves Conceptuelles de résolution des Problèmes Numériques; EDA = Évaluation Des fonctions cognitives et Apprentissages; ERLA = Batterie d'Exploration du Raisonnement et du Langage Associé; EVAC = Épreuves Verbales d'Aptitudes Cognitives; Exalang 3-6 = bilan informatisé pour l'examen du langage et des compétences transversales chez l'enfant de 3 à 6 ans; Exalang 8-11 = Bilan informatisé pour l'examen du langage et des compétences transversales chez l'enfant de 8 à 11 ans; Exalang 11-15 = Batterie informatisée pour l'examen du langage oral, du langage écrit et des compétences transversales chez le collégien; Examath 8-15 = Batterie informatisée d'EXAmen des habiletés MATHématiques; N-EEL = Nouvelles Épreuves d'Examen du Langage; PEDA 1C = Tests pédagogiques de premier cycle primaire; TTR = Tempo Test Rekenen; UDN 2 = Utilisation du nombre, version 2; WIAT-II = Wechsler Individual Achievement Test, 2ème édition; ZAREKI = Die Neuropsychologische Testbatterie für ZAhlenarbeit und REchnen bei Kindern.

Tableau 3

Synthèse des critères et recommandations concernant les outils d'évaluation

Critères		Explications	
Qualification de l'évaluateur		Les qualifications de la personne qui va administrer le test, le corriger et l'interpréter sont clairement explicitées afin de garantir la validité des résultats.	
Standardisation	Consignes de passation et cotation	Les consignes d'administration et de cotation sont clairement spécifiées dans le manuel afin de minimiser la subjectivité lors de l'administration et la cotation.	
	De surface	L'outil est recevable par les utilisateurs. Il s'agit de la compréhension et de l'acceptation du test par les utilisateurs (patients et évaluateurs).	
Validité	De contenu	Validité théorique	La conception de l'outil et le choix des items qui le composent reposent sur les modèles théoriques récents de la fonction cognitive évaluée. Dans le cas d'une batterie de langage comportant quelques subtests mathématiques, la validité de contenu a été établie, non pas pour la batterie au complet, mais pour les subtests en question : 1 point si référence précise et modèle actuel sur la base des données probantes, 0,5 point si référence manquant de précision mais basée sur des données probantes et 0 point si outil s'appuyant sur un modèle théorique non appuyé par la littérature actuelle ou à un outil ne précisant aucune référence.
		Objectif des tests précisé	Les concepteurs posent un choix clair concernant l'objectif de leur outil (diagnostic, détermination d'un niveau de sévérité, orientation thérapeutique) et le précisent.
	De critère	Concomitante	L'outil montre une bonne corrélation entre ses résultats et ceux d'autres épreuves mesurant les mêmes fonctions cognitives et ayant prouvé leur pertinence diagnostique : 1 point est attribué si le test indique les analyses effectuées et si les corrélations indiquées sont moyennes (autour de 0,3) ou bonnes (autour de 0,5); 0,5 point est attribué si les corrélations indiquées sont faibles (c.-à-d. autour de 0,1).
		Prédictive	La pertinence fonctionnelle de l'outil est attestée via une concordance entre les scores observés à l'outil et le fonctionnement dans les activités de vie quotidienne mettant en œuvre la fonction évaluée (p. ex., la note scolaire en mathématique) : 1 point est attribué si le test indique les analyses effectuées et si les corrélations indiquées sont moyennes (autour de 0,3) ou bonnes (autour de 0,5); 0,5 point est attribué si les corrélations indiquées sont faibles (c.-à-d. autour de 0,1).
	De construit	Relations avec les caractéristiques individuelles	Lorsque le construit mesuré est intrinsèquement relié à une ou plusieurs caractéristiques « évidentes » de l'individu, la mesure du construit est sensible à cette relation (sexe, âge, intelligence, etc.) : 1 point est attribué si le test a été évalué selon au moins deux caractéristiques (p. ex., genre, âge, classe, niveau socio-économique, latéralité, pays, présence de trouble ou non); 0,5 point est attribué le test a été évalué selon une seule caractéristique.
		Validité factorielle	Différents items ou sous-tests (malgré des différences de contenu, de format ou de tâches) mesurent une dimension commune qui les influence tous. Le test a la capacité d'établir des associations statistiques entre ses items (ou sous-tests) en conformité avec les dimensions (facteurs) supposément mesurées : 1 point est attribué si une analyse factorielle a été analysé et met en évidence des facteurs

		de regroupement correspondant aux modèles théoriques apportés par les auteurs et aux regroupements en modules effectués par les auteurs.
	Sensibilité/Spécificité	Le pouvoir discriminant de l'outil, c'est-à-dire sa sensibilité et sa spécificité, a fait l'objet d'analyses spécifiques (incluant notamment une population en difficulté), afin de garantir son pouvoir diagnostique. Il s'agit du calcul des pourcentages de vrais positifs, vrais négatifs, faux positifs et faux négatifs : 1 point est attribué si le test indique une sensibilité et une spécificité supérieure à 0,80; 0,5 point si le test indique une sensibilité et une spécificité inférieure à 0,80 est attribué; aucun point n'est attribué si cela n'est pas été testé.
	Temporelle	L'outil fait preuve d'une fidélité test-retest suffisante afin de garantir la stabilité des résultats dans le temps. Friberg (2010) recommande un coefficient de corrélation de 0,90. Un coefficient de 0,80 est acceptable.
	Stabilité	Versions parallèles L'outil montre que deux versions du même test aux mêmes personnes sont équivalentes. L'équivalence d'un test indique à quel point les scores fournis sont indépendants du contenu spécifique des items qui composent le test : 1 point est attribué si le test indique ses analyses et si les corrélations indiquées sont moyennes ou bonnes; 0,5 point est attribué si les corrélations indiquées sont faibles.
	Inter-juges	L'outil fait preuve d'une fidélité inter-juges suffisante afin de garantir que les résultats obtenus sont les plus objectifs possibles et indépendamment de la personne qui a administré et corrigé le test. 1 point a été attribué si les corrélations indiquées étaient égales ou supérieures à 0,90 (recommandation de Friberg, 2010) ou si le Kappa de Cohen indiqué était égal ou supérieur à 0,60 (recommandation de Fleiss, 1981).
Fidélité	Corrélations	Le test montre les corrélations obtenues entre chacun des items du test et le score total au test, ainsi qu'entre les scores totaux de chaque module du test : 1 point est attribué si le test indique ses analyses et si les corrélations indiquées sont moyennes ou bonnes; 0,5 point est attribué si les corrélations indiquées sont faibles.
	Cohérence interne	Bissection L'outil fait preuve de fidélité par bissection, technique qui consiste à diviser un test (une seule version) en deux parties « équivalentes » afin de calculer un « sous-score » pour chacune de ces parties. Les deux parties doivent être corrélées : 1 point est attribué si le test indique ses analyses et si les corrélations indiquées sont moyennes ou bonnes; 0,5 point est attribué si les corrélations indiquées sont faibles).
	Cohérence inter-items	Le manuel fait état d'une analyse statistique de la pertinence des items inclus dans les épreuves, notamment en démontrant la cohérence interne (alpha de Cronbach). Le seuil minimal d'acceptabilité étant estimé à 0,70.
Normes	Taille de l'échantillon (nombre par groupe)	La taille de l'échantillon d'étalonnage est suffisamment importante : 1 point est attribué si l'échantillon comporte au minimum 100 participants par tranche d'âge/sous-groupe; 0,5 point est attribué si l'échantillon comporte au minimum 80 participants par tranche d'âge/sous-groupe; 0,75 point est attribué si certains groupes sont au-dessus de 100 et certains autres à 80.

Description de l'échantillon	Les caractéristiques géographiques, socioéconomiques, linguistiques, l'âge et le genre de la population de l'échantillon d'étalonnage sont clairement explicités : 1 point est attribué si le manuel précise au moins deux caractéristiques (p. ex., genre, âge, classe, niveau socioéconomique, latéralité, pays, présence de trouble ou non); 0,5 point est attribué si le manuel indique une seule caractéristique.
Représentativité de l'échantillon	Les caractéristiques géographiques, socioéconomiques, linguistiques, l'âge et le genre de la population de l'échantillon d'étalonnage sont représentatives de la population tout venant : 1 point est attribué si l'échantillon est représentatif sur au moins deux caractéristiques; 0,5 point est attribué si l'échantillon est représentatif sur une seule caractéristique.
Mesures de tendance centrale	Les moyennes et écarts-types de l'échantillon d'étalonnage sont mentionnés pour chaque tranche d'âge (et/ou les percentiles si la distribution des scores n'est pas gaussienne).
Intervalle de confiance	L'intervalle de confiance (IC) à 95% est rapporté pour chaque norme calculée. C'est un intervalle de valeurs qui a 95% de chance de contenir la vraie valeur du paramètre estimé.

Tableau 4a

*Relevé des caractéristiques psychométriques des différentes batteries d'évaluation analysées : qualification, standardisation et validité**

Nom du test ou de la batterie	Qualification de l'évaluateur	Standardisation : consignes de passation et cotation	Validité							
			De surface	De contenu		De critère		De construit		
				Validité théorique	Objectif des tests précisés	Concomitante	Prédictive	Relations avec les caractéristiques individuelles	Validité factorielle	Sensibilité/Spécificité
Batteries ou tests mathématiques										
B-LM	1	1	0	0	1	0	0	0	0	0
ECPN	0	0,5	0	1	1	0	0	0,5	0	0
ERLA	1	0,5	0	0	1	0	0	0	0	0
Examath 8-15	1	1	0,5	1	1	0	1	1	0	0,5
MathEval	0	1	1	1	1	0	1	0,5	0	0
Numerical	0	1	0	0,5	1	0,5	0	1	1	0
Protocole du calcul élémentaire	0	0	0	1	1	0	0	0	0	0
Tedi-math	1	1	0	0,5	1	0	0,5	0	0	0
Tedi-math Grands	1	1	0	1	1	0	0,5	0	1	0
TTR	1	1	0	1	1	1	1	1	0	0
UDN 2	1	1	0	0	1	0	0	0	0	0
WIAT-II	1	1	0	0	1	0	0	1	0	0
Zareki-R	1	1	0	1	1	0,5	1	0,5	1	0
Zareki-R (article)	1	1	0	1	1	0	0	1	0	0
Nombre de tests remplissant le critère	10	12	1,5	9	14	2	5	6,5	3	0,5
% de tests remplissant le critère	71%	86%	11%	64%	100%	14%	36%	46%	21%	4%
Batteries générales comportant un ou plusieurs subtests mathématiques										
ECHAS	1	1	0	0	0	0	0	1	0	0
EDA	1	1	0	1	1	0	0	1	0	0
EVAC	1	1	0	0	1	0	1	1	0	0
Exalang 3-6	1	1	0,5	0	1	0,5	0,5	0,5	0	0
Exalang 8-11	1	1	0,5	0	1	0,5	0,5	1	0	0
Exalang 11-15	1	1	0,5	0	1	0,5	0,5	1	0	0
N-EEL	1	1	0	0	1	0	0	0,5	0	0
PEDA 1C	1	1	0	0	0	0	0	0	0	0
Nombre de tests remplissant le critère	7	7	1,5	1	6	1,5	2,5	5	0	0
% de tests remplissant le critère	100%	100%	19%	13%	75%	19%	31%	75%	0%	0%
Tout test										

Nombre de tests remplissant le critère	18	20	3	10	20	3,5	7,5	12,5	3	0,5
% de tests remplissant le critère	82%	91%	14%	45%	91%	16%	34%	57%	14%	2%

Notes. *Les tests sont classés en ordre alphabétique selon leur titre. B-LM = Bilan logico-mathématique, cycle 2; ECHAS = Échelle des Apprentissages Scolaires; ECPN = Épreuves Conceptuelles de résolution des Problèmes Numériques; EDA = Évaluation Des fonctions cognitives et Apprentissages; ERLA = Batterie d'Exploration du Raisonnement et du Langage Associé; EVAC = Épreuves Verbales d'Aptitudes Cognitives; Exalang 3-6 = bilan informatisé pour l'examen du langage et des compétences transversales chez l'enfant de 3 à 6 ans; Exalang 8-11 = Bilan informatisé pour l'examen du langage et des compétences transversales chez l'enfant de 8 à 11 ans; Exalang 11-15 = Batterie informatisée pour l'examen du langage oral, du langage écrit et des compétences transversales chez le collégien; Examath 8-15 = Batterie informatisée d'EXAMen des habiletés MATHématiques; N-EEL = Nouvelles Épreuves d'Examen du Langage; PEDAL 1C = Tests pédagogiques de premier cycle primaire; TTR = Tempo Test Rekenen; UDN 2 = Utilisation du nombre, version 2; WIAT-II = Wechsler Individual Achievement Test, 2ème édition; ZAREKI = Die Neuropsychologische Testbatterie für ZAHlenarbeit und REchnen bei KIndern.

Tableau 4b

*Relevé des caractéristiques psychométriques des différentes batteries d'évaluation analysées : fidélité et normes**

Nom du test ou de la batterie	Fidélité						Normes					Score de qualité psychométrique	
	Stabilité			Cohérence interne			Taille de l'échantillon	Description de l'échantillon	Représentativité de l'échantillon	Mesures de tendance centrale	Intervalle de confiance	Nombre de critères validés	% de critères validés
	Temporelle	Versions parallèles	Inter-juges	Corrélations	Bissection	Cohérence inter-items (alpha de Cronbach)							
Batteries ou tests mathématiques													
B-LM	0	0	0	0	0	0	0	1	0,5	0	0	4,5	21%
ECPN	0	0	0	0	0	0	0	1	0	0	0	4	19%
ERLA	0	0	0	0	0	0	0	0	0	0	0	2,5	12%
Examath 8-15	0,5	0	1	0,75	0	0	0,75	1	1	1	1	14	67%
MathEval	0	0	0	1	0	1	0	0,5	0	1	0	9	43%
Numerical	0	0	0	1	0	1	1	1	0	1	0	10	48%
Protocole du calcul élémentaire	0	0	0	0	0	0	1	1	0	0	0	4	19%
Tedi-math	0	0	0	0	0	1	0	1	0,5	1	1	8,5	40%
Tedi-math Grands	0	0	0	0,75	0	0,5	0	1	1	1	1	10,75	51%
TTR	0	0	0	0	0	0	0	1	1	1	0	10	48%
UDN 2	0	0	0	0	0	0	0	1	0	0	0	4	19%
WIAT-II	0	0	0,5	1	1	0	0	1	1	1	1	10,5	50%
Zareki-R	0	0	0	0,5	0	0	0	1	0,5	1	0	10	48%
Zareki-R (article)	0	0	0	0	0	0	0,5	1	1	1	0	8,5	40%
Nombre de tests remplissant le critère	0,5	0	1,5	5	1	3,5	3,25	12,5	6,5	9	4	/	/
% de tests remplissant le critère	4%	0%	11%	36%	7%	25%	23%	89%	46%	64%	29%	/	/
Batteries générales comportant un ou plusieurs subtests mathématiques													
ECHAS	0	0	0	0	0	0	1	1	0	1	0	6	29%
EDA	0,5	0	0	0	0	0	1	1	1	1	0	9,5	45%
EVAC	0	0	0	0	0	0	1	1	0	1	0	8	38%
Exalang 3-6	0,5	0	0	0	0	0	0	1	1	1	0	8,5	40%
Exalang 8-11	0,5	0	0,5	0,5	0	0	1	1	1	1	1	12	57%
Exalang 11-15	0,5	0	0,5	0	0	0	0,5	1	1	1	1	11	52%
N-EEL	0,5	0	0	0	0	0	1	1	1	1	0	8	38%
PEDA 1C	0	0	0	0	0	0	1	1	0	1	0	5	24%

Nombre de tests remplissant le critère	2,5	0	1	0,5	0	0	6,5	8	5	8	2	/	/
% de tests remplissant le critère	31%	0%	13%	6%	0%	0%	81%	100%	63%	100%	25%	/	/
Tout test													
Nombre de tests remplissant le critère	3	0	2,5	5,5	1	3,5	9,75	20,5	11,5	17	6	/	/
% de tests remplissant le critère	14%	0%	11%	25%	5%	16%	44%	93%	52%	77%	27%	/	/

Notes. *Les tests sont classés en ordre alphabétique selon leur titre. B-LM = Bilan logico-mathématique, cycle 2; ECHAS = Échelle des Apprentissages Scolaires; ECPN = Épreuves Conceptuelles de résolution des Problèmes Numériques; EDA = Évaluation Des fonctions cognitives et Apprentissages; ERLA = Batterie d'Exploration du Raisonnement et du Langage Associé; EVAC = Épreuves Verbales d'Aptitudes Cognitives; Exalang 3-6 = bilan informatisé pour l'examen du langage et des compétences transversales chez l'enfant de 3 à 6 ans; Exalang 8-11 = Bilan informatisé pour l'examen du langage et des compétences transversales chez l'enfant de 8 à 11 ans; Exalang 11-15 = Batterie informatisée pour l'examen du langage oral, du langage écrit et des compétences transversales chez le collégien; Examath 8-15 = Batterie informatisée d'EXAMen des habiletés MATHématiques; N-EEL = Nouvelles Épreuves d'Examen du Langage; PEDALC = Tests pédagogiques de premier cycle primaire; TTR = Tempo Test Rekenen; UDN 2 = Utilisation du nombre, version 2; WIAT-II = Wechsler Individual Achievement Test, 2ème édition; ZAREKI = Die Neuropsychologische Testbatterie für Zahlenarbeit und Rechnen bei Kindern.