



HAL
open science

Benchmark techniques vs. science ouverte : à qui revient la parole ?

Brigitte Bigi

► **To cite this version:**

| Brigitte Bigi. Benchmark techniques vs. science ouverte : à qui revient la parole ?. 2025. <hal-05161531>

HAL Id: hal-05161531

<https://hal.science/hal-05161531v1>

Preprint submitted on 14 Jul 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



ETALAB - Open licence

Benchmark techniques vs. science ouverte : à qui revient la parole ?

L'exemple de SPPAS comme réponse concrète aux limites de l'évaluation numérique.

Ce document propose une prise de position sur les limites des benchmarks techniques dans le traitement automatique de la parole, en prenant le logiciel libre SPPAS comme étude de cas illustrant une approche conforme aux principes FAIR de la science ouverte.

Brigitte Bigi
LPL, CNRS, Aix Marseille Univ

15 juillet 2025

Ce document est mis à disposition selon les termes de la [Licence Ouverte 2.0 \(Etalab\)](#), qui autorise la réutilisation, la modification et la diffusion, y compris à des fins commerciales, à condition d'en mentionner la source.

2 **Mots-clés** : SPPAS, segmentation automatique de la parole, benchmarks, science
3 ouverte, FAIR, évaluation, reproductibilité, alignement forcé.

4

5 **Résumé**

6 Les systèmes de segmentation automatique de la parole reposent sur une série d'étapes
7 fondamentales (normalisation, phonétisation, alignement) souvent dissimulées dans
8 des pipelines opaques. Cette invisibilisation limite fortement le contrôle, la
9 reproductibilité et l'adaptabilité des traitements, en particulier dans les contextes de
10 recherche. Parallèlement, les benchmarks standards utilisés pour évaluer ces outils ne
11 prennent en compte qu'un sous-ensemble des dimensions pertinentes, ignorant les
12 phénomènes propres à la parole spontanée, la transparence du processus, ou encore la
13 possibilité d'adaptation aux besoins d'un projet scientifique.

14 Dans ce contexte, SPPAS adopte une approche radicalement différente. Conçu comme
15 un outil de recherche, il offre une architecture ouverte, modulaire et contrôlable, où
16 chaque étape du traitement est explicite et chaque ressource modifiable. Ce
17 positionnement, conforme aux principes FAIR (*Findable, Accessible, Interoperable,*
18 *Reusable*), permet un usage rigoureux, reproductible et collaboratif. En refusant les
19 évaluations réductrices, SPPAS met en lumière la nécessité de repenser les critères de
20 qualité des outils scientifiques. Ce travail défend une évaluation fondée non sur la seule
21 performance technique, mais sur l'utilisabilité réelle, l'ouverture des ressources, et le
22 respect des principes de la science ouverte.

23 Au-delà de la critique, ce travail propose une réponse partielle mais concrète au
24 problème de l'évaluation : une grille multi-critères combinant ouverture, adaptabilité et
25 utilisabilité, offrant une base initiale pour des comparaisons plus pertinentes que celles
26 fondées sur les seuls scores techniques.

27

1.

28

La segmentation automatique de la parole comme tâche invisibilisée

29

Vue rapide :

- **Contexte** : la segmentation de la parole est une tâche ancienne, omniprésente, mais rarement mise en avant et encore mal outillée pour la recherche en SHS.
- **Problème** : la majorité des outils sont jugés sur des scores numériques issus de benchmarks techniques.
- **Thèse** : ces mesures ne reflètent pas les besoins réels des utilisateurs spécialisés car ces métriques techniques sont insuffisantes pour évaluer l'adéquation d'un outil à un usage scientifique.
- **Position** : SPPAS incarne une alternative conçue dans une logique de transparence, adaptabilité, reproductibilité et service à la communauté scientifique.

30

31 La segmentation automatique de la parole est une tâche ancienne, omniprésente dans
32 de nombreux travaux scientifiques, mais paradoxalement rarement mise au premier
33 plan. Historiquement intégrée à la reconnaissance automatique de la parole (RAP), elle
34 a longtemps été considérée comme une étape « intermédiaire », absorbée dans des
35 chaînes de traitement plus larges. Pourtant, elle constitue un moment critique du
36 traitement de la parole, en particulier pour les chercheurs en phonétique expérimentale,
37 en sciences du langage, ou plus généralement en SHS.

38 Dès les années 1950, les premières expériences de reconnaissance phonémique ont
39 intégré implicitement des mécanismes de segmentation. Cependant, ce n'est qu'avec la
40 généralisation des modèles de Markov cachés (HMM), à partir des années 1970, que l'on
41 a vu émerger des approches systématiques et modulaires d'alignement et de découpage
42 temporel de la parole. Ces travaux ont posé les bases des systèmes d'alignement forcé,
43 devenus essentiels dans des domaines aussi variés que la dictée vocale, la synthèse
44 vocale ou l'analyse prosodique.

45 Sur cette base théorique, plusieurs générations d'outils logiciels ont vu le jour. Dans les
46 années 1990, HTK - Hidden Markov Model Toolkit (Young et al., 2002) s'est imposé

47 comme une référence incontournable : puissant – bien que réservé à des utilisateurs
48 avancés, il a profondément influencé la recherche. En parallèle, CMU Sphinx, développé
49 dès les années 1980 à Carnegie Mellon (Seymore et al., 1998), a marqué les débuts de
50 l’open source dans le domaine, en introduisant des systèmes de reconnaissance vocale
51 à large vocabulaire, indépendants du locuteur. À partir de 1999, le projet Julius
52 (Kawahara et al. 2000), développé par le Speech Research Center au Japon, a proposé
53 une solution rapide, légère, et orientée vers les besoins pratiques. Son intégration dans
54 des environnements Unix, sa compatibilité avec la plupart des distributions Linux et son
55 installation directe via *Homebrew* sur macOS en font encore aujourd’hui l’un des outils
56 les plus faciles à installer pour les utilisateurs non spécialistes. Au début des années
57 2010, c’est Kaldi (Povey et al., 2011) qui a redéfini les standards de la recherche moderne
58 en reconnaissance vocale. Très modulaire, s’appuyant sur des scripts bien structurés et
59 des techniques avancées (notamment les réseaux de neurones profonds), Kaldi est
60 devenu la boîte à outils privilégiée des chercheurs, malgré sa complexité de prise en
61 main.

62 Depuis 2020, une nouvelle génération de modèles fondés sur l’apprentissage auto-
63 supervisé a profondément influencé les approches de traitement du signal vocal.
64 Wav2vec 2.0 (Baeovski et al., 2020) a été l’un des premiers à démontrer l’efficacité de
65 représentations acoustiques apprises directement à partir du signal brut, sans
66 supervision alignée. Il a été suivi de modèles comme HuBERT (Hsu et al., 2021) ou
67 Whisper (Radford et al., 2022), qui poursuivent cette logique avec des architectures de
68 type Transformer. Ces modèles sont désormais intégrés dans de nombreux frameworks
69 open source, notamment Hugging Face Transformers (Wolf et al., 2020), et servent de
70 base à des systèmes complets de transcription ou d’alignement. Des projets comme
71 *WhisperX* ou certaines extensions de *pyannote.audio* utilisent ces représentations pré-
72 entraînées pour proposer des segmentations automatiques fines, exploitables dans des
73 contextes applicatifs. Cette évolution marque une rupture avec les chaînes
74 traditionnelles basées sur HMM ou GMM décrites précédemment, en remplaçant
75 l’architecture classique par des blocs neuronaux entraînés de bout en bout.

76 Ces avancées techniques, bien que significatives, répondent en priorité aux besoins de
77 la reconnaissance vocale à grande échelle — transcription, indexation, commande
78 vocale — plutôt qu’à ceux de l’analyse linguistique fine. Or, dans des disciplines comme
79 la phonétique, la linguistique ou l’étude du langage parlé, la segmentation demeure une
80 tâche centrale et encore imparfaitement outillée. Elle ne se limite pas à aligner des
81 phonèmes, mais exige de prendre en compte la diversité des conventions d’annotation,
82 la variabilité des langues et dialectes, ou encore la richesse des phénomènes
83 interactionnels (disfluences, amorces, rires, etc.). Bien que certains considèrent la
84 segmentation comme une tâche « résolue » ou « dépassée », les réalités de terrain
85 montrent au contraire sa persistance comme problème scientifique non trivial — en
86 particulier dans les corpus spontanés, multi-locuteurs, ou issus de contextes
87 minoritaires. Aucun des systèmes mentionnés ci-dessus, aussi puissant soit-il, ne répond
88 aux besoins concrets de segmentation dans un cadre scientifique exploratoire :
89 visualisation, correction, documentation, adaptation au domaine. Ils sont évalués par et

90 pour des développeurs, pas par les chercheurs en SHS qui utilisent les résultats pour
91 produire de la connaissance. Finalement, les systèmes « s'améliorent », mais les
92 utilisateurs ne sont pas mieux servis.

93 Il est donc nécessaire d'apporter un regard critique à l'évaluation de ces systèmes.
94 Aujourd'hui encore, l'évaluation des outils de segmentation repose presque
95 exclusivement sur des indicateurs chiffrés tels que l'UBPA - *Unit Boundary Positioning*
96 *Accuracy*, qui mesure le pourcentage de frontières automatiques situées dans un
97 intervalle (souvent 20 ms) autour des frontières manuelles. Ces métriques techniques
98 sont devenues des standards implicites, valorisés dans les publications, les dépôts de
99 code et les comparaisons entre systèmes. Des travaux ont montré que même entre
100 annotateurs humains, l'accord plafonne à environ 93–94 % dans cet intervalle, ce qui
101 souligne le caractère ambigu et graduel de nombreux phénomènes phonétiques
102 (Hosom, 2009). Si ces chiffres fournissent une forme de repère, ils sont trop souvent
103 interprétés de manière absolue, sans prise en compte du contexte d'usage. Pourtant, ils
104 reflètent très imparfaitement les besoins réels des utilisateurs spécialisés. Ils mesurent
105 une précision numérique, souvent calculée sur des corpus normés, mais ne disent rien
106 de la capacité de l'outil à s'adapter à des usages concrets, à des corpus réels, à des
107 contextes scientifiques variés. Par exemple, une amélioration de 1 % du taux d'accord
108 (de 92 % à 93 %) peut sembler significative, mais son impact pratique dépend
109 entièrement du type d'utilisation :

- 110 • Pour un phonéticien qui corrige manuellement les alignements dans un outil comme
111 Praat (Boerasma, 2001), cette différence n'a aucun effet tangible : dans les deux cas,
112 il doit vérifier l'ensemble des frontières et corriger manuellement les quelques cas
113 fautifs.
- 114 • Pour une analyse automatique de grandes bases de données vocales (extraction de
115 formants, mesures prosodiques), les erreurs d'alignement sont filtrées
116 statistiquement en aval. Une légère différence de précision en amont n'affecte donc
117 pas nécessairement la qualité finale des analyses, car les cas aberrants sont exclus.

118 L'interprétation brute des scores comme indicateur universel de qualité est donc
119 trompeuse. Elle occulte la diversité des usages, des langues, des conventions, et des
120 publics concernés. En focalisant l'attention sur des chiffres "objectifs", elle contribue à
121 invisibiliser des dimensions pourtant essentielles : l'adaptabilité de l'outil, sa
122 transparence, ou sa capacité à intégrer les contraintes spécifiques d'un corpus.

123 Dans ce contexte, le présent document défend la position suivante : un logiciel ne peut
124 pas être jugé uniquement à l'aune de performances numériques.

 Performant ne veut pas dire utilisable

125

126 Un outil peut produire des résultats techniquement excellents sans pour autant être
127 réellement utilisable par celles et ceux à qui il est destiné. Dans le domaine de la
128 segmentation de la parole, de nombreux systèmes — qu'ils soient historiques ou très
129 récents — reposent sur des technologies de pointe, mais imposent encore des
130 conditions d'usage restrictives : installation complexe, dépendance à un système
131 d'exploitation particulier, recours exclusif à la ligne de commande, ou encore nécessité
132 de configurer manuellement des paramètres sans documentation adaptée. Ces
133 exigences, bien qu'acceptables pour des profils informatiques expérimentés, constituent
134 un obstacle important pour une part significative des chercheurs en sciences du langage,
135 en phonétique ou en linguistique interactionnelle, qui n'ont ni le temps ni la formation
136 technique pour s'approprier de tels outils.

137 Cette réalité, souvent négligée dans les publications techniques ou les dépôts open
138 source, souligne une vérité simple mais trop souvent oubliée : la performance
139 algorithmique ne suffit pas. Concevoir un outil pour la recherche scientifique suppose
140 également de prendre en compte l'ergonomie, la diversité des profils utilisateurs, la
141 simplicité d'usage, et l'accessibilité.

142 SPPAS (Bigi 2015) s'inscrit dans une logique radicalement différente : celle d'une
143 recherche logicielle ouverte, fondée sur la transparence des traitements, la modularité
144 des composants, et l'adaptabilité aux besoins spécifiques des chercheurs en phonétique,
145 linguistique, et plus généralement en sciences du langage. Son objectif n'est pas de
146 produire de meilleurs scores sur des benchmarks standards, mais de mieux répondre aux
147 usages concrets de la recherche scientifique de ces domaines. L'ambition initiale de
148 SPPAS a donc été inverse à celle de nombreux outils d'alignement : concevoir un système
149 complet, reproductible, contrôlable, mais suffisamment accessible pour être utilisé par
150 des non-informaticiens. Cela s'est traduit par des choix structurels précis : la présence
151 d'une interface graphique (dont l'utilisation est optionnelle), l'utilisation de fichiers de
152 configuration lisibles et documentés, ou encore l'export direct vers des formats
153 largement adoptés dans les communautés ciblées (comme les fichiers TextGrid
154 compatibles Praat).

155 Et pourtant, malgré ces précautions, il reste des barrières : certains utilisateurs
156 expriment déjà des réticences face à l'installation préalable de Python — preuve que
157 même un outil pensé pour l'usage réel peut rencontrer des freins inattendus. Ce constat,
158 trop souvent négligé dans les publications techniques rappelle une évidence essentielle
159 qui est que **l'utilisabilité ne se mesure pas en millisecondes d'écart sur une frontière
160 de phonème, mais en obstacles que l'on parvient — ou non — à lever pour permettre
161 une appropriation effective.**

162 Cette posture scientifique — fondée sur la transparence, la reproductibilité et
163 l'utilisabilité — est précisément ce que les benchmarks ne mesurent pas, et ce que ce
164 document se propose de défendre.

165 2.

166 Une chaîne invisible mais toujours présente

Vue rapide :

- **Contexte** : Tout outil de segmentation automatique repose sur une séquence d'étapes fondamentales (normalisation, phonétisation, alignement), qui sont systématiquement présentes dans tous les systèmes existants.
- **Problème** : Ces étapes sont souvent dissimulées ou fusionnées dans des boîtes noires logicielles, empêchant leur contrôle, leur compréhension ou leur adaptation par l'utilisateur final.
- **Thèse** : La séparation explicite et la capacité de modifications de ces étapes est essentielle pour garantir un usage scientifique fiable, reproductible et adapté à la diversité des contextes réels.
- **Position** : SPPAS rend visibles, configurables et partageables chacune des étapes du traitement, en accord avec les principes FAIR et dans une optique de service aux besoins concrets des chercheurs.

167

168 Tout système de segmentation automatique, quelle que soit sa forme ou sa
169 sophistication technique, repose inévitablement sur une même séquence d'opérations
170 fondamentales :

- 171 1. Une normalisation du texte brut, qui transforme les chiffres, dates, abréviations
172 ou ponctuations en formes compatibles avec une conversion en phonèmes ;
- 173 2. Une conversion en phonèmes, à l'aide de ressources linguistiques (dictionnaires,
174 grammaires, réseaux de neurones entraînés) ;
- 175 3. Un alignement temporel qui projette la séquence de phonèmes sur le signal
176 audio en exploitant un modèle (traditionnellement un modèle acoustique).

177 Ces trois étapes sont toujours présentes, même lorsqu'elles sont invisibles à l'utilisateur.
178 Dans la majorité des outils récents, elles sont encapsulées dans un pipeline unique,
179 parfois intégré dans un exécutable ou dans une API distante. Cette invisibilisation du
180 pipeline est d'autant plus problématique que les benchmarks actuels n'évaluent en
181 général que la dernière étape (l'alignement), et encore, selon un critère unique (la
182 distance temporelle entre frontières de phonèmes). Toute erreur introduite plus tôt —
183 mauvaise normalisation, phonétisation inadaptée — est alors invisible dans l'évaluation,
184 bien qu'elle puisse fausser l'ensemble du résultat final.

185 Cette dissimulation des étapes n'est pas neutre puisqu'elle prive l'utilisateur de toute
186 vue sur le traitement, l'empêche d'intervenir à un stade donné, et rend toute tentative
187 de validation scientifique difficile, voire impossible. L'analyse linguistique ou phonétique
188 exige pourtant transparence et traçabilité : savoir comment les données sont
189 transformées, pouvoir modifier les ressources utilisées (par exemple le dictionnaire de
190 prononciation), ou encore identifier quelle étape a introduit une erreur ou une
191 approximation.

192 À ce stade, il est essentiel de distinguer deux niveaux de conception dans les systèmes
193 de segmentation automatique :

- 194 • d'une part, les outils algorithmiques comme *HTK*, *Julius* ou *Kaldi*, qui réalisent
195 l'alignement temporel en s'appuyant sur un modèle acoustique et une procédure
196 d'optimisation (souvent via un algorithme de type Viterbi) ;
- 197 • d'autre part, les « *wrappers* » ou environnements utilisateurs comme *SPPAS*, *MFA*
198 ou *WebMAUS*, qui encapsulent ces moteurs pour en faciliter l'usage. Ces wrappers
199 prennent en charge la gestion des ressources linguistiques (dictionnaires, modèles),
200 la configuration des appels, le format des données, l'organisation des sorties, voire
201 l'intégration dans des pipelines d'analyse plus larges.

202 Cette distinction n'est pas seulement conceptuelle : elle est révélatrice d'un besoin
203 structurel de médiation entre la puissance algorithmique d'un moteur et l'usage réel des
204 utilisateurs. Le simple fait que des dizaines de wrappers aient été développés autour de
205 HTK ou Kaldi témoigne de l'inadéquation persistante de ces moteurs bruts avec les
206 besoins concrets de segmentation, d'annotation ou de correction.

207 Un inventaire publié en 2018 (<https://github.com/pettarin/forced-alignment-tools>)
208 recensait déjà plus de 40 outils de ce type — chacun tentant, à sa manière, de combler
209 le fossé entre les exigences techniques de l'alignement automatique et les réalités de
210 l'usage scientifique. Et ce chiffre n'a cessé d'augmenter depuis. Cette prolifération
211 illustre à quel point la question centrale n'est pas seulement "quel algorithme aligne le
212 mieux", mais "*quel outil rend ce traitement réellement exploitable*".

213 Avec l'essor des modèles de reconnaissance vocale fondés sur des réseaux neuronaux
214 profonds, comme *Wav2Vec 2.0*, *HuBERT* ou *Whisper*, un nouveau paradigme technique
215 s'est imposé : la reconnaissance directe à partir du signal, sans explicitation des étapes
216 intermédiaires. Ces modèles, souvent encapsulés dans des bibliothèques comme
217 Hugging Face ou PyTorch, permettent d'obtenir rapidement des transcriptions à partir
218 d'un simple fichier audio. Ils sont ainsi devenus des composants-clés de nouveaux
219 wrappers comme *WhisperX* ou *pyaudio-align*. Pour les informaticiens, cette évolution
220 représente un gain d'utilisabilité considérable par rapport aux moteurs précédents (HTK,
221 Kaldi), souvent lourds à compiler et à configurer. Mais pour les chercheurs non
222 techniciens, la barrière reste intacte : installation en ligne de commande, dépendances
223 multiples, absence d'interface, sorties peu documentées. La promesse de simplicité
224 reste, dans les faits, largement illusoire. Mais surtout, cette nouvelle génération d'outils
225 hérite d'une limite conceptuelle profonde : elle ne voit pas ce qui n'est pas lexical.

226 Les modèles comme *Whisper* ou *Wav2Vec* sont entraînés pour transcrire efficacement
227 de la parole fluide, standardisée, exempte de « perturbations ». Ces modèles ont
228 démontré d'excellentes performances pour la transcription automatique de la parole —
229 mais cette transcription est conçue dans une perspective lexicale, fluide, épurée. Les
230 représentations auto-supervisées qu'ils exploitent sont entraînés sur des corpus
231 normalisés, dépouillés de tout ce qui les gêne. Dans leur logique, tout ce qui relève des
232 rires, soupirs, chevauchements, hésitations, amorces, disfluences — est traité comme
233 du bruit. Ces événements, pourtant au cœur des analyses en linguistique
234 interactionnelle ou en phonétique conversationnelle (par exemple (Bigi & Meunier,
235 2018)), sont soit ignorés, soit activement filtrés. Le pipeline de ces modèles ne prévoit
236 aucune place pour les formes non lexicales, et n'offre aucun mécanisme natif pour les
237 détecter, les annoter ou même simplement les laisser visibles.

238 Ce silence algorithmique est problématique. Il ne s'agit pas d'un simple "manque de
239 performance" : c'est un biais de conception. Ces outils ne ratent pas les disfluences et
240 autres événements — ils ne les cherchent même pas. Et dans des contextes de recherche
241 où ces phénomènes sont centraux, leur usage conduit à une normalisation forcée de la
242 parole, en rupture avec les objectifs mêmes de l'analyse. Car tous ces phénomènes sont
243 au cœur de l'analyse du langage parlé réel :

- 244 • en phonétique (ex. durée, prosodie, pauses)
- 245 • en interaction (ex. gestion du tour de parole)
- 246 • en sociolinguistique (ex. marques de style ou d'émotion)


247 Ce que ces modèles offrent en précision temporelle, ils le perdent en fidélité
248 phénoménologique. Ils sont conçus pour "entendre des mots" — pas pour écouter la
249 parole humaine dans toute sa richesse. **Un outil qui aligne parfaitement des mots, mais**
250 **ignore ce qui fait la spécificité du parler humain, ne peut pas répondre aux attentes**
251 **des chercheurs en sciences du langage.**

252 Chaque phonéticien travaille avec des objectifs, des corpus et des conventions
253 spécifiques. Il n'existe pas de "standard universel" de l'annotation phonétique : certains
254 projets privilégient une segmentation fine de certains phonèmes, d'autres se satisfont
255 de quelques imprécisions au profit d'une validité globalement correcte. Les conventions,
256 comme les objectifs, varient selon les langues, les domaines, ou les finalités de
257 recherche. Cette diversité n'est pas une faiblesse mais une richesse du champ — à
258 condition qu'elle soit prise en compte par les outils.

259 C'est précisément en raison de cette exigence de contrôle, de précision, et d'adaptation
260 que les annotations ont longtemps été réalisées à la main. Le chercheur devait pouvoir
261 décider des unités pertinentes, des critères de segmentation, ou des cas limites à
262 trancher. Or les systèmes de segmentation automatique actuels, lorsqu'ils imposent des
263 chaînes de traitement opaques ou des conventions fixes, reproduisent l'illusion d'un
264 alignement objectif, au détriment de l'ajustement raisonné au terrain. Un bon outil ne
265 doit donc pas imposer un modèle unique. Il doit, au contraire, fournir à l'utilisateur les
266 moyens d'adapter les ressources, d'explorer les variantes, de documenter ses choix.

267 Cette adaptabilité raisonnée n'est pas un luxe, mais une exigence minimale pour garantir
268 la validité scientifique des annotations produites.

269

 **Idée principale** : La segmentation est toujours un processus en plusieurs étapes. Pour qu'un outil soit réellement scientifique, ces étapes doivent être visibles, compréhensibles et ajustables par l'utilisateur.

270

271 Face à cette opacité généralisée des systèmes de segmentation, qu'elle soit technique
272 ou méthodologique, SPPAS adopte une posture radicalement différente, fondée sur les
273 principes **FAIR** (*Findable, Accessible, Interoperable, Reusable*), aujourd'hui reconnus
274 comme un socle des bonnes pratiques en **science ouverte**. Cette exigence de
275 transparence, bien qu'inscrite depuis 2016 dans le droit français via l'article 30 de la loi
276 pour une République numérique -
277 https://www.legifrance.gouv.fr/jorf/article_jo/JORFARTI000033202841 (accès ouvert
278 par défaut aux résultats de la recherche publique), faisait déjà partie intégrante de la
279 conception de SPPAS dès sa première publication en 2012. L'ouverture, la possibilité de
280 modifier les ressources, et de partager les traitements sont au cœur même de sa
281 conception.

282 Dans SPPAS, chaque étape du processus est clairement identifiée : la normalisation du
283 texte, la conversion graphème-phonème et l'alignement temporel sont séparées,
284 explicites et contrôlables. Les ressources utilisées sont toutes fournies, documentées et
285 modifiables — qu'il s'agisse des lexiques, des dictionnaires de prononciation, des
286 modèles acoustiques ou des règles de traitement. Enfin, chaque traitement est traçable
287 et reproductible : aucun paramètre n'est dissimulé, aucune transformation n'est
288 appliquée sans que l'utilisateur puisse la comprendre, l'ajuster ou la corriger. L'ensemble
289 du code source est publié sous licence libre (AGPL), et les ressources sont elles aussi
290 distribuées sous des licences ouvertes compatibles avec la recherche et la réutilisation.
291 **Ce n'est pas un simple choix de conception — c'est une exigence scientifique.** Dans des
292 disciplines où la parole est objet d'analyse fine, où les conventions varient selon les
293 langues, les corpus ou les objectifs de recherche, un outil "boîte noire" est une
294 contradiction. On ne peut pas étudier ce que l'on ne voit pas.



SPPAS rend visible ce que les autres cachent

295

296 Contrairement à la majorité des wrappers, SPPAS a été pensé dès l'origine comme un
297 outil de recherche, au service des usages concrets. Il ne se contente pas de fournir un
298 alignement, mais offre à l'utilisateur un contrôle complet sur le processus : le
299 dictionnaire phonétique peut être enrichi ou modifié, les variantes inappropriées
300 supprimées, les ressources partagées ou adaptées à des corpus spécifiques (langage

301 enfantin, locuteurs dialectaux, parole pathologique...). La transparence est structurelle :
302 les fichiers intermédiaires sont lisibles, les modules configurables, les étapes séparées et
303 explicitées. L'interface graphique permet une prise en main accessible, mais n'enferme
304 pas l'utilisateur dans un parcours contraint.

305 Ce positionnement reflète une conviction forte : un outil scientifique ne doit pas
306 seulement "donner un résultat", mais s'inscrire dans un dialogue avec l'utilisateur, en
307 respectant les exigences de traçabilité, d'adaptabilité, et de reproductibilité, propres à
308 la recherche.

309 La pertinence de SPPAS tient dans sa capacité à s'ajuster aux pratiques de la recherche.
310 Or, **c'est précisément cette dimension — ce qu'un outil rend possible, et non**
311 **uniquement ce qu'il produit — que les benchmarks actuels laissent dans l'ombre.**

 **En résumé, SPPAS, c'est « FAIR » de la science !**

Tous les outils de segmentation réalisent les mêmes étapes (normalisation, phonétisation, alignement), mais trop souvent de manière opaque. SPPAS les rend visibles, modifiables et traçables — conformément aux principes FAIR — ce qui permet aux chercheurs d'ajuster le traitement à leurs besoins réels. Cette adaptabilité, pourtant essentielle, est ignorée par les benchmarks techniques.

312

3.

313

Les limites fondamentales des benchmarks actuels

Vue rapide :

- *Contexte* : Les benchmarks techniques (UBPA, PER, WER, etc.) sont devenus la référence quasi exclusive pour évaluer les outils de segmentation de la parole.
- *Problème* : Ces benchmarks se focalisent sur des écarts locaux par rapport à des références manuelles, sans jamais prendre en compte la transparence, la reproductibilité ni l'adéquation à l'usage réel.
- *Thèse* : Un outil peut obtenir un excellent score tout en restant inutilisable pour les chercheurs en linguistique, phonétique ou SHS.
- *Position* : SPPAS propose une approche radicalement différente, fondée non sur l'optimisation d'un score unique, mais sur la lisibilité et la pertinence de chaque étape pour les usages scientifiques réels.

314

315 Malgré leur technicité apparente, la plupart des systèmes de segmentation automatique
316 reposent sur les mêmes fondements : une chaîne de traitement implicite rarement
317 exposée et jamais évaluée dans sa globalité. Cette invisibilité méthodologique est
318 renforcée par la place dominante qu'occupent aujourd'hui **les benchmarks techniques**,
319 lesquels condensent la qualité d'un outil en quelques **scores numériques standardisés**.
320 Pourtant, ces évaluations masquent des dimensions essentielles : la transparence des
321 traitements, leur traçabilité, leur reproductibilité, ou encore leur adéquation aux
322 données réelles. Cette section examine les limites structurelles des benchmarks actuels,
323 en analysant successivement les trois étapes fondamentales d'un système de
324 segmentation : la normalisation, la phonétisation et l'alignement forcé. Pour chacune
325 d'elles, il s'agit de préciser ce que les benchmarks mesurent — et surtout ce qu'ils ne
326 mesurent pas.

327

3.1 Normaliser pour phonétiser : une étape (trop) invisible

328 Avant toute conversion d'un texte en phonèmes, une étape préalable et indispensable
329 s'impose : la normalisation du texte. Il s'agit d'un ensemble d'opérations visant à
330 transformer un texte brut en une séquence linguistiquement et phonétiquement
331 exploitable. Cette normalisation peut inclure la segmentation en unités, la suppression
332 ou la transformation de signes de ponctuation, le traitement des chiffres, sigles, dates,
333 abréviations, formes contractées ou régionales. Certains outils récents, notamment dans

334 les chaînes Python, appliquent même des transformations particulièrement agressives,
335 comme la romanisation systématique des caractères, ce qui revient à supprimer les
336 accents ou distinctions diacritiques importantes : “côte”, “côté” deviennent
337 indistinctement “cote”, tout comme “élevé” et “élève” deviennent “eleve” — avec des
338 conséquences évidentes sur la phonétisation.

339 Dans les langues où les conventions orthographiques sont stables — comme l’anglais
340 écrit standard — cette étape est souvent intégrée discrètement dans le pipeline, voire
341 passée sous silence. Mais dans la plupart des langues, cette normalisation est loin d’être
342 triviale : en français, les contractions, élisions et variations typographiques dépendent
343 du domaine, du registre ou du style ; dans les langues non segmentées (comme le thaï,
344 le khmer, le chinois traditionnel ou simplifié, ...), l’absence d’espaces entre les mots exige
345 des algorithmes spécifiques de segmentation morphosyntaxique. Dans tous ces cas, la
346 qualité de la normalisation conditionne directement la fiabilité des étapes suivantes.

347 Pourtant, cette étape est rarement évaluée en tant que telle. Elle est soit implicitement
348 incluse dans l’évaluation des modules G2P (grapheme-to-phoneme), soit absente des
349 benchmarks. Les rares études existantes se concentrent souvent sur des cas limités
350 (comme la conversion de nombres en lettres en synthèse vocale), et ne tiennent pas
351 compte des phénomènes présents dans les corpus oraux : abréviations, hésitations,
352 mots inventés, formes incomplètes, etc. De plus, les outils actuels n’offrent que peu ou
353 pas de transparence sur ce qu’ils transforment, ni comment ils le font.

354 SPPAS, au contraire, fait de la normalisation une étape autonome, documentée, publiée
355 (Bigi, 2014) et modulaire du pipeline. Il propose une solution basée sur des modules
356 génériques, valables pour de nombreuses langues ; et des modules spécialisés par
357 langue, remplaçables ou ajustables selon les besoins.

358 Ce choix technique et scientifique permet de trouver un équilibre entre deux extrêmes :
359 d’un côté, une approche “universelle” simpliste qui ignore les spécificités linguistiques ;
360 de l’autre, une approche par langue totalement indépendante, difficile à maintenir et
361 mutualiser. En segmentant clairement cette étape, en en fournissant les composants, et
362 en laissant la possibilité de les modifier ou les remplacer, SPPAS garantit la
363 reproductibilité, l’adaptabilité et la traçabilité du traitement, en accord avec les principes
364 FAIR.

👉 Aucun benchmark actuel ne permet de mesurer la qualité, la modularité ou la pertinence de cette normalisation.

365

366 **3.2 Phonétisation : ce que mesure (réellement) PER**

367 La phonétisation (ou conversion grapheme-to-phoneme, G2P) est l’étape qui consiste à
368 transformer une séquence de mots écrits en une séquence de phonèmes, en vue de
369 l’aligner avec un signal audio ou de le synthétiser. Cette étape repose historiquement :

- 370 • soit sur des dictionnaires phonétiques manuels,
371 • soit sur des règles de prononciation explicites,
372 • soit sur des modèles statistiques ou neuronaux apprenant ces correspondances à
373 partir de données annotées.

374 Si cette tâche paraît bien définie sur des données écrites, lexicales, et standardisées, elle
375 devient beaucoup plus délicate dès que l'on aborde la parole réelle, marquée par des
376 formes inédites, des mots hors lexique, des réductions phonétiques, ou des effets
377 dialectaux.

378 Pour la segmentation automatique de la parole, la phonétisation s'intercale entre la
379 normalisation du texte et l'alignement avec le signal audio, et constitue un maillon
380 indispensable de tout système de traitement de la parole. Pourtant, dans de nombreux
381 outils, cette étape est intégrée de manière invisible : l'utilisateur n'a aucun moyen de
382 savoir comment une transcription phonétique a été produite, ni pourquoi, ni comment
383 la corriger ou l'adapter.

384 Les évaluations classiques de G2P reposent sur des métriques comme le PER - *Phoneme*
385 *Error Rate*, calculé entre une référence et la sortie du système sur des corpus alignés et
386 normés comme CMUdict ou Librispeech (Panayotov et al., 2015). Ces benchmarks
387 mesurent essentiellement la capacité à reproduire une prononciation canonique pour
388 des mots déjà connus. Dans ce cadre, les performances peuvent être impressionnantes.
389 Certains modèles atteignent une PER inférieure à 0,5 % pour des langues comme le
390 coréen, et les architectures de type Transformer surpassent désormais les LSTM sur la
391 plupart des jeux de données standards. Les LSTM - *Long Short-Term Memory*, sont des
392 réseaux de neurones récurrents conçus pour modéliser des séquences, en conservant
393 une mémoire des éléments passés. Ils ont longtemps été l'architecture dominante pour
394 les tâches séquentielles en Traitement Automatique des Langues et en Reconnaissance
395 Automatique de la Parole. Les Transformers (Vaswani et al., 2017) reposent entièrement
396 sur des mécanismes d'attention, sans récursivité. Ils permettent un parallélisme plus
397 efficace et capturent les dépendances longues plus facilement. Depuis 2018, ils ont
398 largement remplacé les LSTM dans la plupart des tâches et commencent à dominer aussi
399 dans la reconnaissance automatique de la parole, la traduction et le G2P.

400 Mais ces scores élevés masquent une réalité : **les benchmarks actuels ne testent que**
401 **des cas idéaux**. Ils ne disent rien sur la capacité à gérer des mots inconnus, des
402 réductions orales, des amorces, des hésitations, ou des variantes dialectales. En bref, ils
403 ne disent rien sur la parole telle qu'elle est effectivement produite et étudiée dans les
404 recherches en linguistique. Or ces formes non canoniques sont particulièrement riches
405 en informations phonétiques, sociolinguistiques ou interactionnelles.

406 Face à cela, SPPAS adopte une approche radicalement différente. Tout d'abord, la
407 phonétisation repose sur un dictionnaire ouvert et modifiable. Pour chaque mot connu,
408 toutes les variantes sont conservées. Si un mot est inconnu, un mécanisme de *longest*
409 *matching* est utilisé pour le découper en segments phonétisables, à partir d'unités
410 présentes dans le dictionnaire (lettres, syllabes, morphèmes). Cela permet à la fois de

411 générer automatiquement des transcriptions plausibles dans les langues peu dotées, et
412 d’atteindre une large couverture lexicale dans les langues mieux pourvues, sans créer de
413 “trou” dans les données.

414 Crucialement, SPPAS ne choisit pas à la place de l’utilisateur. Il n’impose pas une
415 prononciation unique, ni ne cache les variantes. Le choix de la séquence phonétique la
416 plus compatible est délégué à l’algorithme d’alignement, qui exploite le signal pour
417 résoudre l’ambiguïté. Ce découplage entre génération et sélection garantit non
418 seulement une transparence maximale du traitement, mais permet aussi d’obtenir une
419 phonétisation aussi proche que possible de la réalité acoustique observée — sans
420 imposer de standardisation a priori. Là où les benchmarks mesurent l’adhésion à un
421 standard illusoire, SPPAS privilégie l’adaptation au corpus, la liberté de modifier les
422 ressources, et la reproductibilité des résultats. Ce sont ces propriétés, pourtant
423 essentielles pour la recherche, que les benchmarks échouent à mesurer. Enfin, comme
424 pour la normalisation, SPPAS fait de la phonétisation une étape autonome, documentée,
425 publiée (Bigi, 2016) et modulaire du pipeline.

👉 Les benchmarks de G2P évaluent une capacité à produire des séquences de phonèmes corrects sur des mots connus, dans une langue canonique. Ils ignorent tout ce qui caractérise la parole réelle : spontanée, désordonnée, idiosyncratique, dialectale, émotionnelle, inachevée.

426

427 **3.3 Alignement forcé : ce que mesure (réellement) l’UBPA**

428 L’alignement forcé est la tâche qui consiste à associer chaque segment d’un signal audio
429 à une unité linguistique (souvent des phonèmes), en déterminant les bornes temporelles
430 précises de ces unités. Ce processus intervient après la génération d’une ou plusieurs
431 séquences phonétiques issues de la transcription textuelle, par un module de
432 phonétisation. Dans la plupart des systèmes, l’alignement ne sélectionne pas une
433 séquence, mais prend en entrée une grammaire, c’est-à-dire un ensemble de variantes
434 phonétiques possibles, et détermine la séquence la plus probable en fonction du signal.

435 Dans ce contexte, l’UBPA - *Unit Boundary Positioning Accuracy*, est une métrique
436 souvent utilisée pour évaluer la précision temporelle de l’alignement. Elle calcule la
437 proportion de frontières détectées automatiquement qui tombent dans une certaine
438 tolérance (souvent ± 20 ms) autour des frontières définies manuellement. En apparence
439 simple et objective, cette mesure repose sur plusieurs hypothèses implicites qui en
440 limitent fortement la portée. D’abord, pour que l’UBPA soit calculable, il faut fournir une
441 séquence unique de phonèmes correspondant exactement à l’annotation manuelle.
442 Cela neutralise toute variabilité liée au G2P (qui pourrait proposer plusieurs séquences
443 possibles) et réduit artificiellement la tâche à un problème d’alignement pur, dans des
444 conditions rarement réalistes. L’UBPA ne mesure donc ni la capacité du système à choisir
445 la bonne séquence de phonèmes, ni l’adaptabilité à des prononciations non canoniques,

446 ni la gestion des phénomènes oraux typiques (amorphes, chevauchements,
447 interjections...). Ensuite, l'UBPA évalue uniquement le comportement du moteur
448 acoustique et de son modèle statistique, historiquement basé sur des HMM et une
449 recherche Viterbi. Dans les systèmes classiques, cela signifie que le score reflète la
450 qualité du couple *moteur + modèle*, et non celle de l'outil complet. Lorsque les wrappers
451 tels que SPPAS ou WebMAUS encapsulent des moteurs comme HTK ou Julius, c'est ce
452 couple qui est évalué. Le reste de l'environnement — traitements en amont, ergonomie,
453 documentation, interfaces de correction — reste entièrement hors du champ de la
454 métrique. Par exemple, si l'on pouvait extraire un modèle acoustique depuis WebMAUS
455 et qu'on l'utilise dans SPPAS avec le même moteur (HTK en l'occurrence), on obtiendrait
456 des résultats quasi-identiques. L'UBPA ne mesure que les frontières obtenues des
457 phonèmes et ces dernières ne dépendent que du couple *moteur+modèle*.

458 Au-delà de ses limites méthodologiques internes, l'UBPA reste aveugle à de nombreux
459 critères qui conditionnent pourtant l'usage scientifique réel d'un outil de segmentation.
460 Elle ne dit rien de la gestion des phénomènes linguistiques typiques de la parole
461 spontanée — disfluences, amorces, chevauchements, pauses, interjections ou rires —
462 qui sont souvent absents des benchmarks, mais omniprésents dans les corpus de
463 recherche. Elle ignore également la variabilité dialectale, les réductions phonétiques ou
464 les variations stylistiques, autant d'éléments pourtant cruciaux en phonétique,
465 sociolinguistique ou analyse interactionnelle. Surtout, l'UBPA ne prend en compte ni
466 l'adaptabilité d'un outil à des corpus spécifiques, ni la transparence du traitement, ni la
467 reproductibilité des résultats. Elle ne renseigne ni sur la documentation fournie, ni sur la
468 facilité à éditer/adapter les ressources linguistiques, ni sur la capacité à corriger ou
469 visualiser les résultats. Ces dimensions — interface utilisateur, contrôle des paramètres,
470 édition des données — sont pourtant fondamentales dans une démarche scientifique
471 rigoureuse.

472 Pire encore, en valorisant exclusivement la précision temporelle d'un moteur encapsulé,
473 l'UBPA peut involontairement favoriser des systèmes opaques ou rigides, difficiles à
474 adapter, mais optimisés pour un score de benchmark donné. Autrement dit, les outils
475 les mieux classés selon cette métrique sont parfois les moins exploitables pour un
476 chercheur qui souhaite analyser, corriger ou réutiliser les sorties. Surtout, il faut le
477 rappeler avec fermeté : l'UBPA ne mesure en aucun cas la qualité d'un wrapper. Elle
478 évalue uniquement le comportement du couple moteur d'alignement + modèle
479 acoustique. Ainsi, toute étude prétendant comparer des wrappers sur la base de cette
480 seule métrique repose sur une erreur méthodologique : elle confond la performance
481 d'un composant interne avec la valeur scientifique d'un environnement logiciel dans son
482 ensemble.

👉 L'UBPA permet d'évaluer un algorithme d'alignement forcé sur corpus standardisé. Elle ne dit rien sur la transparence du traitement, la modularité des ressources, ou l'adaptabilité aux corpus réels.

483

4.

484

485 Conclusion

486

Vue rapide :

Position : Le respect de la science ouverte et des principes FAIR implique de rendre chaque étape visible, contrôlable, reproductible. C'est cette transparence que les benchmarks ignorent. Dans la recherche scientifique, un bon score ne suffit pas si on ne sait pas comment on l'a obtenu.

487

488 Tout système de segmentation automatique, quelle que soit sa forme ou sa
489 sophistication, repose inévitablement sur une même séquence d'opérations
490 fondamentales : d'abord une normalisation du texte brut, ensuite une conversion en
491 phonèmes, puis un alignement avec le signal audio. Ces étapes sont toujours présentes,
492 même si elles sont souvent dissimulées dans des interfaces simplifiées ou des pipelines
493 opaques. Cette dissimulation n'est pas anodine : elle prive l'utilisateur de toute capacité
494 de contrôle, d'adaptation ou de vérification. Elle transforme une tâche scientifique en
495 une opération algorithmique non reproductible.

496 **Cette opacité des systèmes de traitement de la parole est confortée par des**
497 **benchmarks techniques** — qui ne prennent en compte ni l'ensemble des phénomènes
498 de la parole, ni la transparence, la reproductibilité ou la traçabilité du processus de
499 segmentation. Un outil peut obtenir un "bon score" tout en restant totalement impropre
500 à toute forme d'analyse linguistique ou phonétique. Le véritable choix pour les
501 utilisateurs n'est donc pas entre des outils « plus ou moins performants » au sens des
502 benchmarks, mais entre des systèmes ouverts, adaptables et reproductibles, et des
503 systèmes fermés, rigides, difficilement exploitables en contexte réel.

504 Plusieurs cas pratiques illustrent cette inadéquation. Lors d'une correction manuelle, par
505 exemple, une différence de 6 % ou 7 % sur un score de précision n'a pratiquement aucun
506 effet réel. Dans les deux cas, l'utilisateur est contraint de relire l'intégralité du fichier
507 pour s'assurer de la qualité des frontières. Le gain mesuré ne se traduit donc ni en gain
508 de temps, ni en réduction de l'effort cognitif requis. De même, dans les analyses de
509 grande ampleur, l'extraction des informations comme les valeurs de formants des
510 voyelles est automatisée, puis les valeurs aberrantes sont repérées à l'aide de mesures
511 statistiques afin d'être éliminées. Ainsi, une frontière mal placée lors de l'alignement
512 entraîne l'élimination d'un segment : que l'UBPA soit à 93 % ou à 97 % n'induit aucun
513 changement dans l'analyse qui s'ensuit. Enfin, les corpus de parole spontanée, dialoguée

514 ou pathologique, par exemple, sont traversés de phénomènes que ces outils ignorent
515 systématiquement : réductions, amorces, hésitations, mots inventés, néologismes,
516 expressions dialectales. Or ce sont précisément ces formes atypiques qui intéressent la
517 recherche linguistique, sociolinguistique ou interactionnelle, et que les pipelines
518 opaques tendent à supprimer ou à normaliser.

519 Dans ce contexte, SPPAS propose une alternative radicalement différente. Conçu comme
520 un outil de recherche plutôt qu'un produit logiciel, SPPAS repose sur des principes clairs,
521 rendant chaque étape du processus visible, chaque ressource modifiable, chaque sortie
522 traçable. Ce positionnement, aligné avec les principes **FAIR**, fait de SPPAS un outil adapté
523 à une utilisation scientifique rigoureuse. Cette approche dépasse largement les critères
524 habituels des benchmarks.

 ***Ce que SPPAS fait, les autres ne le font pas*** : SPPAS rend visible, modifiable et partageable chaque étape du processus de segmentation.

525

526 SPPAS a été conçu comme une infrastructure logicielle ouverte, modulaire et
527 contrôlable, en rupture avec les approches opaques des pipelines standards. Il repose
528 sur une séparation explicite entre le moteur d'alignement, interchangeable, et les
529 ressources linguistiques, entièrement accessibles à l'utilisateur. Chaque étape du
530 traitement (normalisation, phonétisation, alignement) a fait l'objet de validations
531 scientifiques documentées et publiées, attestant de la rigueur méthodologique du
532 système. Les ressources utilisées, telles que les dictionnaires phonétiques ou les modèles
533 acoustiques, sont non seulement visibles, mais aussi éditables, remplaçables et
534 partageables. Cette architecture explicite permet une maîtrise fine de l'ensemble du
535 processus, à rebours des systèmes fermés où les choix linguistiques sont dissimulés ou
536 imposés.

537 Ce positionnement a été reconnu au niveau institutionnel. En 2022, SPPAS a reçu un
538 accessit au prix spécial du jury dans le cadre du **Prix Science Ouverte des Logiciels Libres**
539 **de la recherche**, décerné par le Ministère chargé de l'Enseignement Supérieur et de la
540 Recherche. Car le logiciel joue un rôle clé dans la recherche scientifique, dont il est à la
541 fois un outil, un résultat et un objet d'étude. Comme le rappelle le deuxième Plan
542 national pour la science ouverte, « *la mise à disposition des codes sources des logiciels,*
543 *avec la possibilité de les modifier, les réutiliser et les diffuser, est un enjeu majeur pour*
544 *permettre la reproductibilité des résultats scientifiques et soutenir le partage et la*
545 *création de connaissances, dans une logique de science ouverte.* »

546 SPPAS échappe volontairement aux logiques classiques d'évaluation standardisée. Son
547 architecture modulaire, ouverte et configurable rend impossible toute évaluation
548 unidimensionnelle, fondée sur un unique score ou une métrique isolée. Sa richesse
549 fonctionnelle, pensée pour des usages scientifiques réels, déborde largement les cas
550 d'usage figés ou artificiels des benchmarks habituels. Et c'est précisément ce que ces

551 benchmarks échouent à mesurer — la transparence, l’adaptabilité, la traçabilité — qui
552 constitue le cœur de son utilité pour les chercheurs.

553 Il est urgent de dépasser l’hégémonie des benchmarks techniques comme seuls
554 instruments de validation des outils scientifiques. Bien qu’utiles dans certains contextes
555 normés, ces évaluations unidimensionnelles échouent à rendre compte des critères
556 fondamentaux de la pratique de recherche : compréhension des processus, adaptabilité
557 aux données réelles, reproductibilité des résultats. Ce constat appelle à un changement
558 de paradigme dans l’évaluation des outils logiciels de et pour la recherche. Les critères
559 de qualité doivent inclure, aux côtés des performances techniques, des dimensions
560 essentielles telles que :

- 561 • l’ouverture du code et des ressources, condition première d’un usage scientifique
562 transparent ;
- 563 • l’explicabilité des traitements, qui permet à l’utilisateur de comprendre, corriger ou
564 ajuster chaque étape du processus ;
- 565 • la reproductibilité des analyses, rendue possible par la traçabilité complète des
566 données, des paramètres et des algorithmes.

567 Ce sont ces qualités, et non les seuls chiffres de performance sur corpus standardisés,
568 qui garantissent l’adéquation d’un outil aux exigences de la recherche. Valoriser des
569 systèmes capables de s’adapter à la diversité des terrains, de rendre visibles les choix
570 méthodologiques et de respecter les principes de la science ouverte est une nécessité.

571 **Repenser les modalités d’évaluation, c’est aussi affirmer que l’enjeu n’est pas**
572 **seulement technique : il est épistémologique et éthique.**

5.

Perspectives : vers une approche qualitative de l'évaluation

Ce document ne se contente pas de critiquer les limites des évaluations actuelles : il propose une voie concrète pour les dépasser. Il ne s'agit pas seulement de dénoncer l'inadéquation des benchmarks dans certains contextes, mais d'ouvrir un espace de réflexion sur ce que pourrait (et devrait) être une évaluation utile à la recherche.

C'est dans cet esprit que nous proposons, à titre indicatif, un tableau synthétique comparant plusieurs wrappers de segmentation automatique sur des critères élargis. Ces critères ne prétendent pas à l'exhaustivité, ni à l'universalité, mais reflètent des dimensions concrètes, fréquemment rencontrées dans les pratiques scientifiques : transparence, ouverture, contrôlabilité, documentation, respect des licences, possibilités de personnalisation, etc. L'ensemble des critères ont été réparties en 3 dimensions distinctes : 1/ les licences et la conformité légale afin de valider la pertinence avec les critères FAIR, 2/ les fonctionnalités scientifiques afin d'évaluer la pertinence et l'adaptabilité à la tâche, et 3/ l'intégration et la prise en main pour les utilisateurs.

Légende des tableaux :

- ✓ : critère rempli
- ⚠ : critère partiellement rempli ou douteux
- ✗ : critère non rempli
- 🚫 : critère hautement problématique
- ? : non connu/non défini
- 🏆 : prix/distinction
- ⚙ : infrastructure

589

5.1 Licences et conformité

Avant toute considération technique, l'évaluation d'un outil doit porter sur sa légalité d'usage et sa conformité aux principes de la science ouverte. Licence du code, réutilisabilité des ressources linguistiques, respect des licences tierces : ces éléments conditionnent la capacité à diffuser, modifier ou publier des travaux reposant sur l'outil. Ce premier tableau synthétise ces aspects essentiels, souvent absents des publications techniques.

	SPPAS	EasyAlign	WebMAUS	MFA	TorchAudio	Gentle
Licence du code	AGPL v3	?	?	MIT	BSD-2- -Clause	MIT
Licence des ressources linguistiques	GPLv3/ CC	?	?	MIT	MIT/ varié	MIT
Respect des licences tierces	✓	⊘	✓	✓	✓	✓
Code source : F – Findable A – Accessible I – Interoperable R – Reusable	✓ ✓ ✓ ✓	✗ ⚠ ✗ ✗	✗ ✗ ✗ ✗	✓ ✓ ✗ ✓	✓ ✓ ✓ ✓	✓ ✓ ⚠ ⚠
Ressources : F – Findable A – Accessible I – Interoperable R – Reusable	✓ ✓ ✓ ✓	⚠ ✓ ✗ ✗ ✗	⚠ ⚠ ⚠ ✗	✓ ✓ ✗ ⚠	⚠ ✓ ✓ ⚠	✗ ✗ ✗ ✗
Maintenu	✓	✗	✓	✓	✓	✓
Reconnaissance institutionnelle	🏆 MESR	✗	⚙️ CLARIN	✗	✗	✗

597 **Tableau 1** : Statut légal et conformité éthique des logiciels et des ressources

598

599 5.2 Fonctionnalités scientifiques

600 Un bon score à un benchmark ne garantit en rien l'utilité d'un outil pour la recherche.
601 Ce tableau se concentre sur les fonctionnalités linguistiques réellement nécessaires en
602 analyse de la parole : séparation explicite des étapes, gestion des formes atypiques,
603 correction possible, C'est dans ces dimensions que se jouent l'adéquation au terrain,
604 l'adaptabilité, et la transparence du traitement.

	SPPAS	EasyAlign	WebMAUS	MFA	TorchAudio	Gentle
Séparation claire des étapes	✓	✓	✗	✓	✗	✗
Segmentation semi-automatique possible	✓	✗	✗	✗	✗	✗
Phonétisation des mots inconnus	✓	✗	✗	✓	✗	✗
Gestion des variantes de prononciation	✓	✗	✗	✓	✗	✗
Alignement des rires	✓	✗	✗	✗	✗	✗
Gestion spécifique des hésitations	✓	✗	✗	✗	✗	✗

605 **Tableau 2** : Capacités linguistiques et traitement des cas réels

606

607 5.3 Intégration et prise en main

608 Un outil de traitement de la parole destiné à la recherche ne peut se limiter à un moteur
609 performant. Il doit offrir des moyens concrets d'interaction adaptés à divers profils
610 d'utilisateurs : chercheurs sans compétences en programmation, développeurs,
611 ingénieurs. Ce tableau examine les modalités d'usage.

	SPPAS	EasyAlign	WebMAUS	MFA	TorchAudio	Gentle
Interface graphique (web ou locale)	✓	✓	✓	✗	✗	✓
Utilisation via API	✓	✗	✗	✓	✓	✗
Utilisation via CLI	✓	✗	✗	✓	✓	✓
Formats de sortie lisibles	✓	✓	✓	✓	✓	✓

	SPPAS	EasyAlign	WebMAUS	MFA	TorchAudio	Gentle
Personnalisation facile (édition, remplacement)	✓	✓	✗	✓	✓	✗
Documentation développeur	✓	✗	✗	✓	✓	✗
Documentation utilisateur	✓	✓	✓	✓	✗	✗
Tutoriels didactiques	✓	✗	✗	✓	✗	✗
Compatibilité multi-plateforme	✓	✗	✓	✓	✓	✓

612 **Tableau 3** : Souplesse d'usage et modalités concrètes d'interaction

613

614 **⚠ Avertissement** : malgré tous les soins apportés à la collecte et à la vérification
615 des informations, ces tableaux peuvent comporter des erreurs ou des
616 approximations. Ils sont fournis à titre purement indicatif, dans une perspective
617 comparative et non normative.

616

617 L'état des outils ayant pu évoluer, ces données reflètent la situation observée en
618 **juillet 2025** et ne sauraient être tenues pour définitives.

618

619

620 **Références bibliographiques**

- 621 Baevski, A., Zhou, H., Mohamed, A., Auli, M. (2020). wav2vec 2.0: a framework for self-
622 supervised learning of speech representations. Proceedings of the 34th International
623 Conference on Neural Information Processing Systems, Article No.: 1044, pp 12449–
624 12460.
- 625 Bigi, B. (2014). A multilingual text normalization approach. *Human Language Technology*
626 *Challenges for Computer Science and Linguistics, LNAI 8387*, pp. 515–526.
- 627 Bigi, B. (2015). SPPAS-multi-lingual approaches to the automatic annotation of speech.
628 *The Phonetician. Journal of the International Society of Phonetic Sciences, 111*(ISSN:
629 0741-6164), 54-69.
- 630 Bigi, B. (2016). A phonetization approach for the forced-alignment task in SPPAS.
631 *Human Language Technology. Challenges for Computer Science and Linguistics, LNAI*
632 *9561*, pp. 515–526.
- 633 Bigi, B., & Meunier, C. (2018). Automatic segmentation of spontaneous speech. *Revista*
634 *de Estudos da Linguagem, 26*(4).
- 635 Boersma, P. (2001). Praat, a system for doing phonetics by computer. *Glott. Int.*, 5(9),
636 341-345.
- 637 Hosom, J. P. (2009). Speaker-independent phoneme alignment using transition-
638 dependent states. *Speech communication, 51*(4), 352-368.
- 639 Hsu, W. N., Bolte, B., Tsai, Y. H. H., Lakhotia, K., Salakhutdinov, R., & Mohamed, A. (2021).
640 Hubert: Self-supervised speech representation learning by masked prediction of hidden
641 units. *IEEE/ACM transactions on audio, speech, and language processing, 29*, 3451-3460.
- 642 Kawahara, T., Lee, A., Kobayashi, T., Takeda, K., Minematsu, N., Sagayama, S., ... &
643 Shikano, K. (2000). Free software toolkit for Japanese large vocabulary continuous
644 speech recognition. In Proc. Int'l Conf. on Spoken Language Processing (ICSLP), Vol. 4,
645 pp. 476--479.
- 646 Panayotov, V., Chen, G., Povey, D., & Khudanpur, S. (2015). *Librispeech: An ASR corpus*
647 *based on public domain audio books*. In *2015 IEEE International Conference on Acoustics,*
648 *Speech and Signal Processing (ICASSP)* (pp. 5206–5210). IEEE.
649 <https://doi.org/10.1109/ICASSP.2015.7178964>
- 650 Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., ... & Vesely, K.
651 (2011). The Kaldi speech recognition toolkit. In *IEEE 2011 workshop on automatic speech*
652 *recognition and understanding*, vol. 1.
- 653 Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2022). *Robust*
654 *speech recognition via large-scale weak supervision*. OpenAI.
655 <https://openai.com/research/whisper>

656 Seymore, K., Rosenfeld, R., Chen, S., Eskenazi, M., Gouvea, E., Reddy, R., ... & Thayer, E.
657 (1998). The 1997 CMU Sphinx-3 English broadcast news transcription system.

658 Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T.,
659 Louf, R., Funtowicz, M., & Rush, A. M. (2020). *Transformers: State-of-the-art natural*
660 *language processing*. In *Proceedings of the 2020 Conference on Empirical Methods in*
661 *Natural Language Processing: System Demonstrations* (pp. 38–45). Association for
662 Computational Linguistics. <https://doi.org/10.18653/v1/2020.emnlp-demos.6>

663 Young, S., Evermann, G., Gales, M., Hain, T., Kershaw, D., Liu, X., ... & Woodland, P.
664 (2002). The HTK book. *Cambridge university engineering department*, 3(175), 12.

665 **Liens utiles (juillet 2025) :**

666 **CMUdict** - Carnegie Mellon University. *The CMU Pronouncing Dictionary*. Version 0.7b :
667 <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

668 **CMUSphinx** : <https://cmusphinx.github.io/>

669 **Julius** CSR Engine : <https://github.com/julius-speech/Julius>

670 **Gentle** : <https://github.com/strob/gentle>

671 **HTK** – Hidden Markov Toolkit : <https://htk.eng.cam.ac.uk/>

672 **MFA** - Montreal Forced Aligner : [https://github.com/MontrealCorpusTools/Montreal-](https://github.com/MontrealCorpusTools/Montreal-Forced-Aligner)
673 [Forced-Aligner](https://github.com/MontrealCorpusTools/Montreal-Forced-Aligner)

674 **Praat** : <https://praat.org/>

675 **pyannote.audio** : Neural building blocks for speaker diarization :

676 <https://github.com/pyannote/pyannote-audio>

677 **PyTorch** – TorchAudio – Alignement forcé :

678 [https://docs.pytorch.org/audio/main/generated/torchaudio.functional.forced_align.ht](https://docs.pytorch.org/audio/main/generated/torchaudio.functional.forced_align.html)
679 [ml](https://docs.pytorch.org/audio/main/generated/torchaudio.functional.forced_align.html)

680 [https://github.com/pytorch/audio/blob/main/examples/tutorials/forced_alignment tu](https://github.com/pytorch/audio/blob/main/examples/tutorials/forced_alignment_tutorial.py)
681 [torial.py](https://github.com/pytorch/audio/blob/main/examples/tutorials/forced_alignment_tutorial.py)

682 **SPPAS** : <https://sppas.org/>

683 **Wave2vec** : [https://ai.meta.com/blog/wav2vec-20-learning-the-structure-of-speech-](https://ai.meta.com/blog/wav2vec-20-learning-the-structure-of-speech-from-raw-audio/)
684 [from-raw-audio/](https://ai.meta.com/blog/wav2vec-20-learning-the-structure-of-speech-from-raw-audio/)

685 **WebMAUS** :

686 <https://clarin.phonetik.uni-muenchen.de/BASWebServices/interface/WebMAUSBasic>

687 **WhisperX** - Time-Accurate Speech Transcription with Word-Level Timestamps :

688 <https://github.com/m-bain/whisperx>