



HAL
open science

Transcription de corpus oraux d'arabe parlé en interaction. Convention AraPI et annexes.

Lina Choueiri, Loubna Dimachki, Catherine Pinon, Véronique Traverso

► To cite this version:

Lina Choueiri, Loubna Dimachki, Catherine Pinon, Véronique Traverso. Transcription de corpus oraux d'arabe parlé en interaction. Convention AraPI et annexes.. 2019. hal-02153116

HAL Id: hal-02153116

<https://hal.science/hal-02153116>

Preprint submitted on 17 Jun 2019

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.

Transcription de corpus oraux d'arabe parlé en interaction : La convention AraPI

Auteurs :

Véronique Traverso, Catherine Pinon, Loubna Dimashki, Lina Choueiri

Présentation du projet AraPI

- *L'équipe (groupe AraPI – Arabe Parlé en Interaction) :*
 - Ifpo (Véronique Traverso, Catherine Pinon)
 - Université libanaise (Loubna Dimashki, Moustafa al-Hajj)
 - AUB (Lina Choueiri)
 - CNRS - laboratoire ICAR (Joseph Dichy, Carole Etienne)
 - CNRS - laboratoire SEDYL (Stefano Manfredi)
 - CNRS - laboratoire LLL (Layal Kanaan).
- *Les financements :*
 - Ifpo
 - Labex ASLAN
 - Laboratoire ICAR (CNRS)
 - Université libanaise
- *Les objectifs :*
 - établir une convention de transcription pour les corpus d'arabe parlé en interaction
 - alimenter une base de données (corpus transcrits et alignés)
- *Les publications :*
 - Traverso, Pinon, Dimachki et Choueiri : « Corpus d'arabe parlé (1) : Quels corpus d'arabe parlé en libre accès ? », *les carnets de l'Ifpo, la recherche en train de se faire à l'Institut français du Proche-Orient* [En ligne] <https://ifpo.hypotheses.org/8865>
 - Traverso, Pinon, Dimachki et Choueiri : « Corpus d'arabe parlé (2) : contraintes et problèmes liés à la réalisation des corpus d'arabe parlé en interaction », *les carnets de l'Ifpo, la recherche en train de se faire à l'Institut français du Proche-Orient* [En ligne] <https://ifpo.hypotheses.org/9039>
 - Traverso, Pinon, Dimachki et Choueiri : « Corpus d'arabe parlé (3) : Choix pour la notation des sons et des phénomènes dans la réalisation de corpus d'arabe parlé », à paraître dans *les carnets de l'Ifpo, la recherche en train de se faire à l'Institut français du Proche-Orient* [En ligne] <https://ifpo.hypotheses.org/9305>
 - Pinon, Traverso, Dimachki et Choueiri : « Corpus d'arabe parlé (4) : La convention de transcription ARAPI pour l'arabe parlé en interaction », à paraître dans *les carnets de l'Ifpo, la recherche en train de se faire à l'Institut français du Proche-Orient* [En ligne].

ARABE PARLE EN INTERACTION (ARAPI) CONVENTIONS

Introduction	3
0.....Présentation générale de la convention	4
0. 0. <i>Forme de la transcription</i>	4
0. 1. <i>Une transcription en plusieurs tiers</i>	4
0. 2. <i>Les phénomènes interactionnels</i>	5
0. 3. <i>L'alphabet de transcription</i>	5
0. 3. 1. <i>Consonnes</i>	6
0. 3. 2. <i>Voyelles</i>	6
1. La tier 1 : une transcription interactionnelle	7
1. 1. <i>Phénomènes interactionnels</i>	7
1. 2. <i>Transcription des sons et représentation de la réalité sonore (granularité)</i>	7
1. 3. <i>Notation de la longueur</i>	7
1. 3. 1. <i>Longueur phonologique et longueur expressive</i>	7
1. 3. 2. <i>Les voyelles longues finales de la graphie arabe</i>	7
1. 3. 3. <i>Les alif-s suscrits de l'orthographe arabe</i>	8
1. 4. <i>Les déterminants</i>	8
1. 4. 1. <i>L'article</i>	8
1. 4. 2. <i>Les démonstratifs</i>	8
1. 5. <i>Les pronoms</i>	8
1. 5. 1. <i>Les pronoms personnels</i>	8
1. 5. 2. <i>Le pronom relatif</i>	8
1. 5. 3. <i>Les pronoms interrogatifs</i>	8
1. 5. 4. <i>Les pronoms démonstratifs</i>	9
1. 6. <i>Les particules et prépositions</i>	9
1. 7. <i>Les chiffres</i>	9
1. 8. <i>Autres éléments conventionnels</i>	9
1. 8. 1. <i>Les expressions figées</i>	9
1. 8. 2. <i>Les mots ou syntagmes cités en langue étrangère</i>	9
1. 8. 3. <i>Les mots étrangers entrés dans les usages</i>	9
1. 8. 4. <i>Les emprunts arabisés</i>	9
1. 8. 5. <i>Les régulateurs, marques d'hésitation, acquiescement, négation</i>	10
1. 8. 6. <i>Les noms propres</i>	10
1. 8. 7. <i>La mention de noms anonymisés</i>	10
2. La tier 2 : une transcription morpho-phonologique	11
2. 1. <i>Phénomènes interactionnels</i>	11
2. 2. <i>Transcription des sons</i>	11
2. 3. <i>Notation de la longueur</i>	12
2. 3. 1. <i>Longueur phonologique et longueur expressive</i>	12
2. 3. 2. <i>Les voyelles longues finales de la graphie arabe</i>	12
2. 3. 3. <i>Les alif-s suscrits de l'orthographe arabe</i>	12
2. 4. <i>Les déterminants</i>	12
2. 4. 1. <i>L'article</i>	12
2. 4. 2. <i>Les démonstratifs</i>	12
2. 5. <i>Les pronoms</i>	12
2. 5. 1. <i>Les pronoms personnels</i>	12
2. 5. 2. <i>Le pronom relatif</i>	13
2. 5. 3. <i>Les pronoms interrogatifs</i>	13
2. 5. 4. <i>Les pronoms démonstratifs</i>	13

2. 6. Les particules et prépositions.....	13
2. 7. Les chiffres.....	13
2. 8. La morphologie verbale.....	13
2. 9. La morphologie nominale.....	14
2. 10. Autres éléments conventionnels.....	14
2. 10. 1. Le lexique dialectal.....	14
2. 10. 2. Les expressions figées.....	14
2. 10. 3. Les mots ou syntagmes en langue étrangère.....	14
2. 10. 4. Les mots étrangers entrés dans les usages.....	14
2. 10. 5. Les emprunts arabisés.....	14
2. 10. 6. Les régulateurs, marques d'hésitation, acquiescement, négation.....	15
2. 10. 7. Les noms propres.....	15
2. 10. 8. La mention de noms anonymisés.....	15
3. La tier 3 : une transcription en graphie arabe.....	16
3. 1. Phénomènes interactionnels.....	16
3. 2. Transcription des sons.....	16
3. 3. Les voyelles longues finales de la graphie arabe.....	16
3. 4. Les déterminants.....	16
3. 4. 1. L'article.....	16
3. 4. 2. Les démonstratifs.....	17
3. 5. Les pronoms.....	17
3. 5. 1. Les pronoms personnels.....	17
3. 5. 2. Le pronom relatif.....	17
3. 5. 3. Les pronoms interrogatifs.....	17
3. 5. 4. Les pronoms démonstratifs.....	17
3. 6. Particules et prépositions.....	17
3. 7. Les chiffres.....	17
3. 8. La morphologie verbale.....	18
3. 9. Autres éléments conventionnels.....	18
3. 9. 1. Les expressions figées.....	18
3. 9. 2. Les mots ou syntagmes en langue étrangère.....	18
3. 9. 3. Les mots étrangers entrés dans les usages.....	18
3. 9. 4. Les emprunts arabisés.....	18
3. 9. 5. Les régulateurs, marques d'hésitation, acquiescement, négation.....	18
3. 9. 6. Les noms propres.....	19
3. 9. 7. La mention de noms anonymisés.....	19
4. La tier 4 : une traduction.....	20
4. 1. Phénomènes interactionnels.....	20
4. 2. Vouvoiement et ajustements de traduction.....	20
4. 3. Expressions figées.....	20
4. 4. Mots ou syntagmes en langue étrangère.....	20
4. 5. Régulateurs, marques d'hésitation, acquiescement, négation.....	20
5. Liste des documents annexes.....	21
Annexe 1 : Convention de transcription pour les phénomènes oraux et interactionnels.....	21
Annexe 2 : Liste des notations conventionnelles.....	21
Annexe 3 : Liste des gloses morpho-syntaxiques.....	21
Annexe 4 : Règles générales de transcription des dialectes en caractères arabes.....	21
Annexe 5 : Tableau synoptique de la convention AraPI.....	21
6. Bibliographie.....	21

INTRODUCTION

Les transcriptions de corpus d'interaction en arabe parlé pour des analyses de linguistique interactionnelle ou pragmatique en sont encore à leurs balbutiements. Chaque chercheur est confronté aux mêmes difficultés, et forge ses propres solutions en fonction des caractéristiques de ses données et de ses objectifs. Les choix effectués, eux, sont très variés, s'étageant autour de quatre possibilités majeures : une transcription phonétique, qui cherche à rendre compte de façon précise des spécificités sonores des productions ; une transcription phonologique, qui ne retient que les variantes distinctives de la variété concernée ; une translittération, c'est-à-dire une représentation fidèle, lettre par lettre, de l'orthographe arabe en alphabet latin ; et enfin une transcription en caractères arabes.

La convention ARAPI a opté pour une transcription en quatre tiers pour rendre les corpus transcrits exploitables par des chercheurs de différents domaines, éventuellement non arabophones, y compris en recourant aux outils informatiques sur script latin ou arabe. Les quatre tiers concernent 1) une transcription interactionnelle ; 2) une transcription morpho-phonologique assortie d'une glose ; 3) une transcription en caractères arabes ; 4) une traduction en français. Chacune des tiers est présentée en détail dans le document qui suit.

La convention est établie à partir de corpus syro-libanais, mais elle est conçue pour être étendue à d'autres dialectes arabes, moyennant son adaptation à leurs spécificités. De même, la convention peut être utilisée par les chercheurs de différents horizons, quelle que soit leur spécialité : chaque chercheur transcrit selon ses besoins, en sélectionnant la ou les tiers qui sont nécessaires à sa recherche, voire en ajoutant des éléments conventionnels qu'il doit expliciter, lorsqu'un phénomène qu'il a besoin de représenter n'est pas prévu dans la convention.

0. PRESENTATION GENERALE DE LA CONVENTION

Ce document détaille les conventions établies pour la transcription de sources produites en arabe parlé en interaction. Il est augmenté de plusieurs annexes. La première détaille la convention de notation pour les phénomènes oraux. La deuxième propose la liste de tous les éléments linguistiques conventionnels sous forme de tableaux. Cette liste est amenée à évoluer avec l'ajout de nouvelles transcriptions au corpus. La troisième annexe présente les règles générales pour la transcription des dialectes, inspirées de CODA (Conventional Orthography for Dialectal Arabic, Habash et al. 2012¹). La quatrième récapitule les gloses morphosyntaxiques utilisées dans la tier 2b. Et enfin, la cinquième est un tableau synoptique qui permet d'appréhender rapidement les différences notables entre les tiers.

0. 0. Forme de la transcription

La police utilisée pour les tiers 1 et 2a est Courier New, 10 pt, en noir. Pour la tier 2b, 3 et 4, on utilise Times New Roman 11. On insère un saut de ligne entre chaque tier ².

Les tiers commencent par l'abréviation du pseudonyme du locuteur (généralement en trois lettres, ex. AMI), suivies d'une tabulation, puis du texte prononcé par le locuteur. Toute parole transcrite doit être rapportée au locuteur qui l'a produite (voir le détail des mises en page dans l'annexe 1).

Exemple :

AMI → transcription interactionnelle (tier 1)
transcription morpho-phonologique (tier 2a)
glose morpho-syntaxique (tier 2b)
كتابة عربية (tier 3)
traduction dans la langue du chercheur (tier 4)

0. 1. Une transcription en plusieurs tiers

Pour l'établissement d'un corpus complet, la transcription est réalisée en quatre tiers.

Tier 1 : transcription interactionnelle

La tier 1 est une transcription adaptée aux analyses pragmatiques, discursives et interactionnelles. Elle fonctionne en lien avec les données primaires (enregistrements), et elle indépendante de la tier 2. La transcription choisie est phonético-phonologique ; elle conserve le découpage en mots et ne prétend pas rendre compte de toutes les variantes de prononciation. Elle applique également certaines conventions de segmentation des morphèmes et indique les emprunts aux langues étrangères.

¹ Habash, Nizar & Diab, Mona & Rambow, Owen. (2012). Conventional Orthography for Dialectal Arabic. Proceedings of the Language Resources and Evaluation Conference (LREC), Istanbul.

² Ces informations sont pertinentes pour un travail sous word. Elles n'ont plus lieu d'être pour des transcriptions effectuées à l'aide de logiciels de transcription, qui disposent d'interfaces à multiples tiers (comme ELAN ou Praat).

Tier 2 dédoublée : transcription morpho-syntaxique

La tier 2 correspond à une représentation phonologique d'un standard du dialecte concerné qui intègre un certain nombre de choix conventionnels (tier 2a). Elle est assortie d'une glose morpho-syntaxique suivant la référence LeipzigGlossingRules (tier 2b). Voir l'annexe 4 pour la liste des étiquettes morpho-syntaxiques et des exemples en arabe libanais.

Tier 3 : transcription en caractères arabes

La tier 3 est une transcription en caractères arabes s'inspirant de la convention CODA (Conventional Orthography for Dialectal Arabic) ³. Voir dans l'annexe 3 les règles générales pour la transcription des dialectes en caractères arabes.

Tier 4 : traduction

Enfin la tier 4 propose une traduction. Elle se veut fidèle au sens davantage qu'à la structure morpho-syntaxique de l'énoncé. Certains éléments difficilement traduisibles font l'objet de notations conventionnelles.

0. 2. Les phénomènes interactionnels

La notation des phénomènes interactionnels ne concerne que la tier 1. Celle-ci comporte tous les éléments relatifs à la production de la langue parlée en interaction : chevauchements de parole, pauses, enchaînements immédiats, notation sommaire de la prosodie.

Pour ces phénomènes, la convention AraPI reprend la convention ICOR, qui s'inspire des conventions de Jefferson (2004) généralement utilisées en linguistique interactionnelle, avec de légères modifications. Ces conventions sont présentées dans l'annexe 1, Conventions de transcription pour les phénomènes oraux et interactionnels.

Une version étendue des conventions interactionnelles est consultable sur le site CORINTE, <http://icar.univ-lyon2.fr/projets/corinte/>.

0. 3. L'alphabet de transcription

La convention ARAPI utilise l'alphabet phonétique international (API) avec quelques modifications :

- les voyelles longues sont notées avec un trait suscrit (ā, ū, ī, etc.) plutôt qu'avec les deux points (a:, u:, i:, etc.) déjà employés pour la transcription interactionnelle ;
- les consonnes emphatiques sont marquées avec des points sous les consonnes non-emphatiques correspondantes (ṣ, ḍ, ṭ, ḏ). La transcription de l'interdentale emphatique devrait être le ḏ avec un point souscrit. Nous la noterons ḏ̣.

Dans les tableaux ci-dessous, pour des raisons pratiques, les sons de la langue sont présentés selon l'ordre de l'alphabet arabe.

³ Habash, Diab et Rambow, "Unified Guidelines and Resources for Arabic Dialect Orthography" [En ligne] <http://www.lrec-conf.org/proceedings/lrec2018/pdf/395.pdf>

Habash, Diab et Rambow, "Conventional Orthography for Dialectal Arabic" [En ligne] http://www.lrec-conf.org/proceedings/lrec2012/pdf/579_Paper.pdf

0. 3. 1. Consonnes

Arabe	API	Exemple	Arabe	API	Exemple
ب	b, p	baba, paspor	ط	ɬ	ɬawla
ت	t	tīn	ظ	ð, ɖ, ʒ	ðohor, ɖəhr, būza
ث	θ, t, s	θaqafe, talʒ, masalan	ع	ʕ	ʕənd
ج	ʒ, dʒ, g	ʒamāl, dʒamāl, gamāl	غ	ɣ	ɣēm
ح	ħ	ħabīb	ف	f, v	ʔalf, narvaz
خ	x	ʔəxt	ق	q, g, ʔ	ħuqūq, grīb, ʔarīb
د	d	dār	ك	k, g	ktāb, garāʒ, sangari
ذ	ð, d, z	hāða, hāda, ʔəstēz	ل	l, ɭ	lēl, aɭlāh
ر	r	rās	م	m	mabrūk
ز	z	zēt	ن	n	naʒāh
س	s	səne	ه	h	hənnen
ش	ʃ	ʃams	و	w	walad
ص	ʂ	ʂēf	ي	j	jōm
ض	ɖ, ʒ	ɖamir, maʒbūɬ	ء	ʔ	suʔāl

Des caractères arabes particuliers peuvent être utilisés pour la transcription en arabe des mots étrangers courants contenant les sons [p] ou [v].

Arabe	API	Exemple
پ	p	پويات (des pots)
ڤ	v	ڤاڤا (ça va)
گ	g	سنگري (plombier)

0. 3. 2. Voyelles

API	Exemple	API	Exemple
a	abjaɖ	ā	ktāb
ɛ	bərke	ē	lēʃ, bēt, bēred
e	ɬajjeb	ī	dīk
ə	sətt, təffēha	ō	ʃōb
i	kīfik	ū	sūʔ
o	bjektob	ǎ	ǎ... ok (exclamation)
u	suʔāl	ō	bōʒūr
y	styɖjo	∅	ʔ∅ (hésitation)

Le son [e] en finale de mot, lorsqu'il ne s'agit pas d'une *imāla* mais d'un terme étranger, peut être transcrit به en arabe. Exemple : kafe / كافيه (café).

1. LA TIER 1 : UNE TRANSCRIPTION INTERACTIONNELLE

La tier 1 est une transcription phonético-phonologique en caractères latins (alphabet phonétique international adapté), dans laquelle les phénomènes oraux et interactionnels sont notés.

1. 1. Phénomènes interactionnels

Tous les éléments relatifs à la production de la langue parlée en interaction sont notés dans cette tier selon la convention détaillée en annexe 1. Les segments inaudibles sont rendus par xx. En cas de troncation, on transcrit phonétiquement ce que l'on entend. Les transcriptions incertaines sont notées entre parenthèses.

1. 2. Transcription des sons et représentation de la réalité sonore (granularité)

Les sons sont notés tels qu'ils sont perçus :

- le *qaf* est noté (q, ʔ ou g), selon sa réalisation.
- les interdentes (θ, ð), sont notées comme elles sont prononcées (s, d, t, z, θ, ð).
- les emphatiques (ṣ, ḍ, ṭ, ḏ) ne sont notées que si la consonne est réellement vélarisée. Sinon, on notera leur réalisation (éventuellement : s, d, t, z).
- le [l] emphatique dans *alḷāh* est noté avec un point.
- la *hamza* initiale n'est pas transcrite ; la *hamza* médiane et finale n'est transcrite que si on la prononce.
- l'*imāla* est notée par le son [e].

La transcription de la tier 1 ne prétend pas atteindre une granularité très fine pour la représentation des sons prononcés, notamment les voyelles (c'est pourquoi on parle de représentation phonético-phonologique plus que réellement phonétique). Néanmoins, il n'y a pas d'essai de standardisation des sons transcrits ; le transcripteur note simplement ce qu'il perçoit.

1. 3. Notation de la longueur

1. 3. 1. Longueur phonologique et longueur expressive

On distingue la longueur phonologique (distinctive) notée par un trait sur la voyelle (ā, ī, ē, ō, ū) de la longueur expressive notée par les deux points après la voyelle (ba:). La longueur expressive n'apparaît que dans la tier 1.

1. 3. 2. Les voyelles longues finales de la graphie arabe

Les voyelles longues en finale de mot dans la graphie arabe ne sont pas transcrites avec une voyelle longue, mais conformément à leur réalisation. Cela concerne notamment :

- les pronoms personnels
- le *alif maqṣura* ou le *alif ṭawīla* à la fin des verbes, des pronoms, des noms et des particules
- la préposition *fi*
- certains pronoms suffixes.

1. 3 3. Les alif-s suscrits de l'orthographe arabe

On note la voyelle longue uniquement si elle est prononcée comme telle. Cela concerne certains démonstratifs, la particule *lākin* et le mot *allāh*. Le nom Allah est transcrit comme il est prononcé.

Ex. *aḷḷa / aḷḷah / aḷḷā / aḷḷāh*

1. 4. Les déterminants

1. 4. 1. L'article

L'article est noté comme il est prononcé et sans séparateur.

Ex. *lbēt*

1. 4. 2. Les démonstratifs

Le démonstratif réduit est noté [ha], accolé à l'article sans séparateur.

Ex. *halbēt*

Les démonstratifs complets sont notés tels qu'ils sont réalisés.

Ex. *hajda / haza ; hajde / hajdi / hazih ; hadōl / hadōle / hajdōle ; hdāk / hadīke ; hāj* (liste non exhaustive).

1. 5. Les pronoms

1. 5. 1. Les pronoms personnels

Les pronoms sont transcrits selon leur réalisation. Les pronoms suffixes sont transcrits tels qu'ils se prononcent et sont accolés au verbe, au nom ou à la préposition.

Les pronoms compléments introduits par la particule [ijjā] sont notés conformément à leur réalisation, collés à la particule, cette dernière étant séparée de ce qui précède par une espace. Exemple de transcription des pronoms (non exhaustif) :

Pronoms isolés		Pronoms possessifs		Pronoms compléments	
ana, ani	hijsse	ktēbe	ktēba	jāni, jāne	jē, ijjē,
enta	nehna, ?ēhna	ktēbek, ktēbak	ktēbna	jāk, ijjāk, jēk	jāha, jēha
ente	ento	ktēbek, ktēbik	ktēbkon	jāki, ijjāki,	jāna, jēna
huwwe	henne, hennen	ktēbo	ktēbon	jēke	jekon
				jē, ijjē, jāho	jon, jāhon

Ex. *bejton, ?āletlo, fīki, beddik jēha*

1. 5. 2. Le pronom relatif

Les relatifs sont notés au plus proche de la perception et sont séparés du mot suivant par une espace.

Ex. *jilli / əlli / li / l*

1. 5. 3. Les pronoms interrogatifs

Les interrogatifs sont notés selon leur prononciation.

Ex. *kīf, wejn, wēn, fēn, lē, lēf, fū*

1. 5. 4. Les pronoms démonstratifs

Les pronoms démonstratifs sont notés tels qu'ils sont réalisés.

Ex. *hēj*, ; *hajda* ; *hajde* / *hajdi* ; *hadōl* / *hadōle* / *hajdōle* ; *haydēk* / *hadīke* (liste non exhaustive).

1. 6. Les particules et prépositions

Les particules sont notées selon leur prononciation. Elles sont séparées du nom par une espace, même si elles sont accolées dans la graphie arabe. Voir la liste des prépositions dans l'annexe 2.

Ex. *bi lbēt*

1. 7. Les chiffres

Les chiffres sont transcrits phonétiquement.

1. 8. Autres éléments conventionnels

Certains éléments qui sont transcrits phonétiquement en tier 1 feront l'objet de notations conventionnelles dans les autres tiers. C'est le cas du lexique dialectal courant sans équivalent en arabe standard.

1. 8. 1. Les expressions figées

Les éléments à traduction variée (comme *jaʕni*) sont transcrits en tier 1 mais font l'objet d'une notation conventionnelle dans les autres tiers. On trouvera en annexe 2 une liste des notations conventionnelles.

1. 8. 2. Les mots ou syntagmes cités en langue étrangère

Les productions en langue étrangère sont délimitées par des dièses afin de permettre leur repérage rapide et systématique. Ils sont notés, sans mention de langue, dans l'orthographe de la langue concernée pour les segments (courts ou longs) dont la prononciation respecte le standard de la langue concernée, ou en transcription phonétique si l'on considère que cela apporte des informations pertinentes ou si la langue n'est pas identifiable.

On utilise, en cas de besoin, l'orthographe adaptée décrite dans la convention ICOR, qui fait intervenir l'antiquote ` pour les élisions non standard (ex. je l` fais). Voir Annexe 1.

1. 8. 3. Les mots étrangers entrés dans les usages

Les mots étrangers entrés dans les usages sont transcrits phonétiquement, sans dièses.

Ex. *ḅāj*

1. 8. 4. Les emprunts arabisés

Les emprunts arabisés sont transcrits phonétiquement et notés entre dièses.

1. 8. 5. Les régulateurs, marques d'hésitation, acquiescement, négation

La notation des régulateurs, de l'hésitation, de l'acquiescement ou de la négation fait l'objet de conventions. Voir Annexe 2.

1. 8. 6. Les noms propres

Qu'ils soient remplacés par des pseudonymes ou non, les noms propres sont transcrits dans leur orthographe courante si elle est connue ou transcrits phonétiquement selon les cas et notés entre guillemets hauts, avec majuscule initiale. Il peut s'agir d'une personne, d'un lieu, d'une institution, etc.

1. 8. 7. La mention de noms anonymisés

Les noms propres anonymisés qui sont cités dans les interactions apparaissent entre guillemets hauts, en entier ou sous leur forme abrégée selon les cas. Cela peut concerner une personne, un lieu, une institution, etc. S'il s'agit du pseudonyme de l'un des participants à l'interaction, on le conserve sans guillemets.

2. LA TIER 2 : UNE TRANSCRIPTION MORPHO-PHONOLOGIQUE

La tier 2a est une transcription morpho-phonologique d'un standard du dialecte concerné. Elle est notée en caractères latins (alphabet phonétique international adapté) et intègre de nombreuses notations conventionnelles. Elle est assortie, en tier 2b, d'une glose inspirée des Leipzig Glossing Rules (voir l'annexe 4 et le texte intégral [En ligne] <https://www.eva.mpg.de/lingua/pdf/Glossing-Rules.pdf>). Elle comporte de nombreuses notations conventionnelles détaillées en annexe 2. Les choix de notations conventionnelles respectent le plus souvent possible la convention de notation des dialectes arabes CODA et se réfèrent à l'ouvrage *Parlons arabe libanais* de Fida Bizri (L'Harmattan, 2010).

2. 1. Phénomènes interactionnels

Les phénomènes interactionnels ne sont pas notés, à l'exception des pauses notées par (.) quelle que soit leur durée et des chevauchements marqués par [.

Les segments inaudibles sont rendus par xx.

En cas de troncation, on ajoute un tiret à la fin du mot inachevé.

Les transcriptions incertaines sont notées entre parenthèses.

2. 2. Transcription des sons

La tier 2 est un équivalent en caractères latins de la tier 3, elle-même établie en suivant la convention de notation des dialectes arabes CODA (voir l'annexe 3 et la présentation de la tier 3 ci-dessous). Cette tier repose sur de nombreuses notations conventionnelles, puisqu'elle correspond, dès que possible, à une translittération du *ductus* de l'arabe standard, tout en préservant de manière conventionnelle les caractéristiques du dialecte dans lequel l'interaction a lieu.

Les consonnes sont transcrites conformément au *ductus* standard et les voyelles sont transcrites phonologiquement. De ce fait, dès que le mot dialectal a son équivalent en arabe standard, on le translittère selon l'orthographe classique, et lorsque ce n'est pas le cas, on le transcrit phonologiquement selon une notation conventionnelle renvoyant à un standard dialectal. Les notations conventionnelles concernent notamment les pronoms personnels, les particules, les interrogatifs, les démonstratifs, les désinences verbales, etc. (voir l'annexe 2).

Voici les règles générales de transcription des sons :

- Le *qaf* est noté q.
- Les interdentes seront notées θ ou ð.
- Les emphatiques ṣ, ḍ, ṭ, ḏ.
- La *hamza* est notée ʔ en médiane et finale de mot uniquement, et non notée à l'initiale.
- Les sons vocaliques sont réduits aux sons [a] [u] [i] [e] [o] pour les voyelles brèves, à [ā] [ū] [ī] [ē] [ō] pour les longues.
- L'*imāla* est transcrite par le son [e].
- Les diphtongues sont rétablies.
- Les voyelles élidées sont rétablies.
- Les voyelles épenthétiques sont supprimées.

2. 3. Notation de la longueur

2. 3. 1. Longueur phonologique et longueur expressive

Seule la longueur phonologique est notée en tier 2.

2. 3. 2. Les voyelles longues finales de la graphie arabe

Les voyelles longues finales de l'orthographe arabe standard, même si elles ne sont pas prononcées, sont notées dans la tier 2 qui privilégie une translittération de l'orthographe arabe standard.

Ex. *fī* ; *ilā* ; *ʕalā* ; *maqḥā*

2. 3. 3. Les alif-s suscrits de l'orthographe arabe

On note systématiquement une voyelle longue si celle-ci correspond à un *alif* suscrit dans l'orthographe standard. Cela concerne certains démonstratifs, la particule *lākin* et le mot *allāh*.

Ex. *lākin*, *hāḏā*, *allāh*

2. 4. Les déterminants

2. 4. 1. L'article

L'article est segmenté au moyen d'un tiret. En début d'énoncé, il est systématiquement noté [el-]. En cas de liaison, il est toujours noté [-l-] même en cas d'assimilation.

Ex. *el-bēt* ; *bi-l-bēt* ; *bi-l-rās*

2. 4. 2. Les démonstratifs

Le démonstratif *ha* est noté avec séparateur (ha-l-).

Ex. *ha-l-kitāb* ; *bi-ha-l-bajt*

Les démonstratifs complets font l'objet d'une notation conventionnelle.

Ex. *hāḏā* ; *hāḏihi* ; *haḏūl* ; *haḏāk* ; *haḏīk*

2. 5. Les pronoms

2. 5. 1. Les pronoms personnels

Les pronoms suffixes sont segmentés par un tiret.

Les pronoms isolés et suffixes apparaissent sous une forme dialectale conventionnelle.

Les pronoms compléments introduits par la particule [ijjā] seront notés de manière conventionnelle sous la forme *jē*-pronom conventionnel.

Pronoms isolés		Pronoms possessifs		Pronoms compléments	
anā	hija	kitāb- i	kitāb- hā	jē- nī	jē- hā
enta	nehnā	kitāb- ka	kitāb- nā	jē- ka	jē- nā
enti	entō	kitāb- ki	kitāb- kon	jē- ki	jē- kon
huwa	hennen	kitāb- hu	kitāb- hon	jē- hu	jē- hon

2. 5. 2. Le pronom relatif

Le relatif est systématiquement noté *jalli*, séparé du mot suivant par une espace, quelle que soit sa réalisation.

2. 5. 3. Les pronoms interrogatifs

Les interrogatifs font l'objet d'une transcription conventionnelle, quelle que soit leur réalisation. Voir Annexe 2.

2. 5. 4. Les pronoms démonstratifs

Les pronoms démonstratifs font l'objet d'une notation conventionnelle.

Ex. *hāḏā ; hāḏī ; haḏūl ; haḏāk ; haḏīk*

2. 6. Les particules et prépositions

Un tiret sépare les prépositions monolithères du mot qui suit. Elles sont notées conventionnellement (*la-*, *bi-*, *wa-*, *ta-*) ; les autres prépositions sont séparées du nom par une simple espace. Voir la liste des prépositions dans l'annexe 2.

2. 7. Les chiffres

Les chiffres sont transcrits conformément à leur équivalent en arabe standard.

Ex. *arbaṣat ṣafar*

2. 8. La morphologie verbale

Le radical du verbe est isolé de l'ensemble des affixes par un tiret. On segmente aussi les morphèmes de modes/temps et personnes qui font l'objet d'une notation conventionnelle (voir annexe 2).

Ex. *b-ja-ktub-ū*

Les préfixes de conjugaison comme *ṣam* et *raḥ* sont séparés du verbe par une espace.

Ex. *ṣam ja-ktub ; raḥ ta-ktub*

Le morphème *ḥa-* est segmenté par un tiret du préfixe de conjugaison.

Ex. *ḥa-ja-ktub*

En cas d'élision du morphème de conjugaison en tier 1, on rétablit en tier 2 le morphème manquant.

Ex. *ṣam fakkir* (tier 1) = *ṣam u-fakkir* (tier 2a)

De même, pour l'inaccompli à valeur modale (sans le préverbe b-) employé après différentes particules comme حتى, لازم, بدّ, في, etc., la *hamza* de la première personne du singulier est rétablie dans la tier 2 (*badd-ī aḥṣuz*).

La voyelle du radical et celle des personnes pour les formes dérivées sont notées conformément à la voyelle présente dans le verbe correspondant en standard (voir annexe 2).

Ex. *b-ja-ʕmal* ; *b-ju-ballif*

Pour le détail des conventions concernant la transcription des voyelles de conjugaison et des verbes conjugués à l’accompli, à l’inaccompli et à l’impératif, voir l’annexe 2.

2. 9. La morphologie nominale

Les marques du genre et du nombre des noms sont segmentées par un tiret.

Ex. *ʕāmil* ; *ʕāml-e* ; *ʕāml-īn*

2. 10. Autres éléments conventionnels

Certains éléments qui sont transcrits phonétiquement en tier 1 font l’objet de notations conventionnelles dans les autres tiers. Voir Annexe 2.

2. 10. 1. Le lexique dialectal

Les productions typiquement dialectales font l’objet de notation conventionnelle.

2. 10. 2. Les expressions figées

Pour les éléments très récurrents à l’oral à traduction variée (comme les particules discursives, par exemple *jaʕni*) ou les expressions figées d’inspiration coranique ou non, on utilise une notation conventionnelle composée de quatre lettres capitales (voir annexe 2).

Ex: *b-fūf-ki* INCH

2. 10. 3. Les mots ou syntagmes en langue étrangère

Toutes les productions en langue étrangère sont notées entre dièses afin de permettre leur repérage rapide et systématique. Elles sont notées sans mention de la langue, dans l’orthographe de la langue concernée, que les segments soient courts ou longs, ou en transcription phonétique si l’on ignore la langue en question.

Ex. ENNO Nicolas #hair# badd-hu #bonnet#

2. 10. 4. Les mots étrangers entrés dans les usages

Les mots étrangers entrés dans les usages, dont on fournit une liste non exhaustive en annexe 2, seront notés entre dièses, dans l’orthographe de la langue et sans mention de la langue.

Ex. #bonjour# ; #sorry# ; #merci#

2. 10. 5. Les emprunts arabisés

Les emprunts arabisés, dont on fournit une liste non exhaustive en annexe 2, seront notés entre dièses en transcription et sans mention de la langue.

Ex. #pojēt# ; #taperwarāt#

2. 10. 6. Les régulateurs, marques d'hésitation, acquiescement, négation

La notation des régulateurs, de l'hésitation, de l'acquiescement ou de la négation font l'objet de conventions, présentées en annexe 2.

2. 10. 7. Les noms propres

Qu'ils soient remplacés par des pseudonymes ou non, les noms propres sont transcrits dans leur orthographe courante si elle est connue et notés entre guillemets hauts, avec majuscule initiale. Il peut s'agir d'une personne, d'un lieu, d'une institution, etc.

2. 10. 8. La mention de noms anonymisés

Les noms propres anonymisés qui sont cités dans les interactions apparaissent entre guillemets hauts, en entier (avec une majuscule) ou sous leur forme abrégée selon les cas. Cela peut concerner une personne, un lieu, une institution, etc. S'il s'agit du pseudonyme de l'un des participants à l'interaction, on le conserve sans guillemets.

3. LA TIER 3 : UNE TRANSCRIPTION EN GRAPHIE ARABE

La tier 3 propose une transcription en caractères arabes. Elle reprend la convention CODA pour la transcription des dialectes, avec quelques aménagements (voir l'annexe 3). Selon cette convention, on se réfère, chaque fois que cela est possible, au *ductus* de l'arabe standard. En même temps, on conserve les traits spécifiques du dialecte utilisé, en faisant émerger un standard dialectal par le biais de différentes notations conventionnelles.

3. 1. Phénomènes interactionnels

Les phénomènes interactionnels ne sont pas notés, à l'exception des pauses notées par (.) quelle que soit leur durée, et des chevauchements marqués par]. Les segments inaudibles sont rendus par xx. En cas de troncation, on ajoute un tiret à la fin du mot inachevé. Les transcriptions incertaines sont notées entre parenthèses. Si le nom des intervenants est cité dans un énoncé, il figure en caractères latins en trois lettres majuscules, selon le code donné par le transcripteur.

3. 2. Transcription des sons

La graphie standard est rétablie dès que cela est possible : on réintègre toutes les consonnes selon l'orthographe classique.

- Le *qaf* est noté ق, quelle que soit sa réalisation.
- Les interdentales sont notées ث et ذ, quelle que soit leur réalisation.
- Les emphatiques sont notées ص, ض, ط, ظ.
- La *hamza* initiale n'est notée que si elle correspond à une *hamza* stable.
- Les géminations morphologiques sont systématiquement notées par une *chadda*.
- La transcription n'est pas vocalisée (ni voyelles brèves, ni *alif* suscrit).
- Aucun signe de lecture n'est ajouté pour l'article (*hamza, waṣṣla, sukūn* et *chadda*).
- Le *tanwīn* [an] est noté.

Pour les mots qui n'ont pas de graphie standard et certains termes dialectaux courants, on se réfère à une notation conventionnelle, présentée dans l'annexe 2.

3. 3. Les voyelles longues finales de la graphie arabe

Les voyelles longues finales figurent conformément au *ductus* classique. Cela concerne notamment :

- les pronoms personnels dans leur notation conventionnelle (أنا)
- le *alif maqṣūra* ou le *alif ṭawīla* à la fin des verbes, des pronoms, des noms et des particules (على)
- certains pronoms suffixes (ها).

3. 4. Les déterminants

3. 4. 1. L'article

L'article est toujours noté اـ, sans aucun signe de lecture, conformément à la graphie standard.

3. 4. 2. Les démonstratifs

Le démonstratif *ha-* est noté accolé à l'article هـالـهـ.

Les autres démonstratifs font l'objet d'une notation conventionnelle, quelle que soit leur réalisation. Voir l'annexe 2.

3. 5. Les pronoms

3. 5. 1. Les pronoms personnels

Les pronoms sont notés de manière conventionnelle.

Pronoms compléments		Pronoms possessifs		Pronoms isolés	
إياها	إياني	كتابها	كتابي	هي	أنا
إيانا	إياك	كتابنا	كتابك	نحننا	إنت
إياكن	إياكي	كتابكن	كتابكي	إنتوا	إنتي
إياهن	إياه	كتابهن	كتابه	هنن	هو

3. 5. 2. Le pronom relatif

Le relatif sera systématiquement noté اللي.

3. 5. 3. Les pronoms interrogatifs

Les interrogatifs font l'objet d'une transcription conventionnelle, quelle que soit leur réalisation :

كيف، وين، فين، مين، ليه، ليش، إيش، شو، قديش، قديه، إيمنى، إمتى...

Voir l'annexe 2 pour une liste complète.

3. 5. 4. Les pronoms démonstratifs

Les pronoms démonstratifs font l'objet d'une notation conventionnelle, quelle que soit leur réalisation. Voir l'annexe 2.

هذا، هذي، هذول، هذاك، هذيك

3. 6. Particules et prépositions

Les particules et prépositions monolitères sont attachées conformément à l'orthographe arabe.

Les particules tronquées ne sont pas rétablies et sont traitées comme des monolitères.

Ex. : على البيت عالبيت .

3. 7. Les chiffres

Les chiffres sont transcrits conformément à leur équivalent en arabe standard.

Ex. . أربعة عشر

3. 8. La morphologie verbale

La transcription des verbes conserve les marqueurs modaux du dialecte (*b-*, *ʕam*). En cas d'élision du morphème de conjugaison en tier 1 (*ʕam fakkir*), le morphème manquant est rétabli (عم أفكر). De même, pour l'inaccompli à valeur modale (sans le préverbe *b-*) employé après différentes particules comme حتى لازم، بدّ، في، etc., la *hamza* de la première personne du singulier est rétablie dans la tier 3 (بدّي أحجز).

Les tableaux complets de conjugaison à l'accompli, à l'inaccompli et à l'impératif figurent dans l'annexe 2.

3. 9. Autres éléments conventionnels

Des conventions orthographiques sont introduites pour le vocabulaire dialectal ainsi que pour les emprunts (Voir annexe 2).

3. 9. 1. Les expressions figées

Les éléments à traduction variée ou les expressions figées d'inspiration coranique ou non, sont écrits en toutes lettres au plus près de l'orthographe standard, selon une notation conventionnelle qu'on s'attachera à respecter pour faciliter les recherches automatiques (Annexe 2).

3. 9. 2. Les mots ou syntagmes en langue étrangère

Les productions en langue étrangère sont notées avec les caractères et dans l'orthographe de la langue concernée.

Ex. #bonnet# بدّه #hair# إنّه نقولاس

3. 9. 3. Les mots étrangers entrés dans les usages

Les mots étrangers entrés dans les usages seront transcrits en caractères arabes, entre dièses, sans menton de la langue.

Ex. #فاسا# ، #مرسی# ، #سوري# ، #بنجور#

3. 9. 4. Les emprunts arabisés

Les emprunts arabisés seront transcrits en caractères arabes, entre dièses, sans menton de la langue.

Ex. #تپیروارات# ، #پیویات#

3. 9. 5. Les régulateurs, marques d'hésitation, acquiescement, négation

La notation des régulateurs, de l'hésitation, de l'acquiescement ou de la négation font aussi l'objet de conventions. Voir Annexe 2.

3. 9. 6. Les noms propres

Qu'ils soient remplacés par des pseudonymes ou non, les noms propres sont transcrits dans leur orthographe courante si elle est connue et notés entre guillemets hauts, en caractères arabes. Il peut s'agir d'une personne, d'un lieu, d'une institution, etc.

3. 9. 7. La mention de noms anonymisés

Les noms propres anonymisés cités apparaissent en entier ou sous leur forme abrégée selon les cas, entre guillemets hauts et en caractères arabes. Cela peut concerner une personne, un lieu, une institution, etc. S'il s'agit du pseudo de l'un des participants à l'interaction, on le conserve en caractères latins et sans guillemets.

4. LA TIER 4 : UNE TRADUCTION

La tier de traduction offre un accès direct au sens. Les choix de traduction doivent être guidés par deux principes : l'intelligibilité et la conformité à l'original. On ne fera donc pas du mot-à-mot et on veillera à ce que le rendu soit compréhensible en français, tout en s'assurant de rester au plus proche de la structure arabe originale. Par ailleurs, aucun signe de ponctuation n'est ajouté.

4. 1. Phénomènes interactionnels

Les phénomènes interactionnels ne sont pas indiqués, à l'exception des pauses notées par (.) quelle que soit leur durée, et des chevauchements marqués par [. Les segments inaudibles sont rendus par xx. En cas de troncation, on ajoute un tiret à la fin du mot inachevé. Les éléments incertains dans la transcription sont traduits entre parenthèses. Comme dans les autres tiers, si le nom des intervenants est cité dans un énoncé, il figure en caractères latins en trois lettres majuscules, selon le code donné par le transcripteur.

4. 2. Vouvoiement et ajustements de traduction

Les énoncés où le vouvoiement serait d'usage en français sont traduits en utilisant "vous", plutôt que "tu".

Les ajouts aidant à comprendre le sens de la traduction seront mentionnés entre {}.

4. 3. Expressions figées

Les éléments à traduction variée (comme *jaʕni*) ou les expressions figées d'inspiration coranique ou non, sont notés conformément à la notation conventionnelle établie utilisée en tier 2, composée de quatre lettres capitales (voir la liste en annexe 2).

Ex : LAYK je passe te voir demain INCH

4. 4. Mots ou syntagmes en langue étrangère

Les mots prononcés dans une langue autre que l'arabe sont traduits en français, mais figurent entre # pour marquer l'emprunt, la langue utilisée étant indiquée entre double parenthèses à l'intérieur des dièses en se référant au code ISO. De même, les réalisations en arabe classique (CLA) ou standard moderne (MSA) figurent entre dièses avec mention de la langue entre double parenthèses.

Ex : #merci ((FR))#

4. 5. Régulateurs, marques d'hésitation, acquiescement, négation

La notation des régulateurs, de l'hésitation, de l'acquiescement ou de la négation fait l'objet de conventions (voir l'annexe 2).

4. 6. Les noms propres

Qu'ils soient remplacés par des pseudonymes ou non, les noms propres sont notés dans leur orthographe courante si elle est connue, sans guillemets mais avec majuscules initiales.

5. LISTE DES DOCUMENTS ANNEXES

Annexe 1 : Convention de transcription pour les phénomènes oraux et interactionnels

Annexe 2 : Liste des notations conventionnelles

Annexe 3 : Liste des gloses morpho-syntaxiques

Annexe 4 : Règles générales de transcription des dialectes en caractères arabes

Annexe 5 : Tableau synoptique de la convention AraPI

6. BIBLIOGRAPHIE

- BICKEL, B., COMRIE, B., HASPELMATH, M. (2015). *Leipzig Glossing Rules* [En ligne] <https://www.eva.mpg.de/lingua/pdf/Glossing-Rules.pdf>
- BIZRI, F. (2010). *Parlons l'arabe libanais*, Paris, L'Harmattan.
- BLANC, H. (1960). "Stylistic Variations in Spoken Arabic: A Sample of Interdialectal Educated Conversation", in C. A. FERGUSON (ed.), *Contributions to Arabic Linguistics*, Cambridge, Massachusset, Harvard Middle Eastern Monographs, p. 81-161.
- CAUBET, D. (1999). « Arabe maghrébin : passage à l'écrit et institutions », *Faits de Langue* 13, p. 235-244.
- DIAB, M., HABASH, N. (2014). *Natural Language Processing of Arabic and its Dialects* [En ligne] http://emnlp2014.org/tutorials/4_notes.pdf.
- DICHY, J. (1994). « La pluriglossie de l'arabe », *Bulletin d'études orientales* 46, Damas, p. 19-42.
- HABASH, N., SOUDI, A., BUCKWALTER, T. (2007). "On Arabic Transliteration", in SOUDI A., BOSCH A., NEUMANN G. (eds). *Arabic Computational Morphology. Text, Speech and Language Technology*, 38. Springer, Dordrecht.
- HABASH, N., DIAB, M., RAMBOW, O. (2012). *Conventional orthography for dialectal Arabic* [En ligne] http://www.lrec-conf.org/proceedings/lrec2012/pdf/579_Paper.pdf.
- HABASH, F., et al. (2018). *Unified Guidelines and Resources for Arabic Dialect Orthography* [En ligne] <http://www.lrec-conf.org/proceedings/lrec2018/pdf/395.pdf>.
- ICOR (2013). *Convention ICOR* [En ligne] http://icar.cnrs.fr/projets/corinte/documents/2013_Conv_ICOR_250313.pdf
- IZRE'EL, S., METTOUCHI, A. (2015). "Representation of speech in Corpafras. Transcriptional strategies and prosodic units", in METTOUCHI A., VANHOVE M., CAUBET D. (dir.), *Corpus-based Studies of Lesser-described Languages: The CorpAfroAs corpus of spoken AfroAsiatic languages*, John Benjamins Publishing Company, p. 14-41.

- JEFFERSON, G. (2004). "Glossary of transcript symbols with an introduction", in LERNER G. (dir.), *Conversation Analysis: Studies from the First Generation*, Amsterdam, John Benjamins, p. 13-31.
- MAAMOURI, M, GRAFF, D., JIN, H., CIERI, C. (2004a). *Dialectal Arabic Orthography-based Transcription & CTS Levantine Arabic Collection*. EARS PI Meeting and RT-04 Workshop, IBM Executive Conference Center, Palisades, NY, USA; November 7-11, 2004 [En ligne] [http:// www.sainc.com/richtrans2004/](http://www.sainc.com/richtrans2004/)
- MAAMOURI, M, GRAFF, D., JIN, H., CIERI, C. (2004b). "Dialectal Arabic Telephone Speech Corpus: Principles, Tool design and Transcription Convention", in Proceedings of the NEMLAR Arabic Language Resources and Tools Conference Sept. 22-23, 2004. Cairo, Egypt, p. 55-60 [En ligne]
https://www.researchgate.net/publication/228876943_Dialectal_arabic_telephone_speech_corpus_Principles_tool_design_and_transcription_conventions
- METTOUCHI, A., VANHOVE, M., CAUBET, D. (dir.) (2015). *Corpus-based Studies of Lesser-described Languages: The CorpAfroAs corpus of spoken AfroAsiatic languages*. John Benjamins Publishing Company.
- METTOUCHI, A., CHANARD, C. (2010). "From Fieldwork to Annotated Corpora: the CorpAfroAs Project", *Faits de Langue-Les Cahiers 2*, p. 255-265.
- TRAVERSO, V. (2006). *Des échanges ordinaires à Damas*. Lyon/Damas: PUL/Presses de l'IFPO.
- VAILLANT P., LEGLISE I. (2014). « À la croisée des langues. Annotation et fouille de corpus plurilingues », *Revue des nouvelles technologies de l'information*, Cepaduès-Editions, RNTI-SHS-2, p. 81-100.
- ZAGHOUBANI W. (2014). "Critical Survey of the Freely Available Arabic Corpora", Workshop on Free/Open-Source Arabic Corpora and Corpora Processing Tools [En ligne] <https://arxiv.org/ftp/arxiv/papers/1702/1702.07835.pdf>
- ZAWAYDEH Bushra et al. "Guidelines for Transcribing Arabic Dialects" [En ligne] <https://catalog.ldc.upenn.edu/docs/LDC2005S08/BBN-Babylon-transcription-guidelines.pdf>
- ZRIBI I. , GRAJA, M., ELLOUZE, M., JAOUA, M., HADRICH BELGUITH, L., (2013). "Orthographic Transcription for Spoken Tunisian Arabic", *CICLing 1*, p. 153-163.
- ZRIBI I., BOUJELBANE R., MASMOUDI A., ELLOUZE M., HADRICH BELGUITH L., HABASH N. (2014). "A Conventional Orthography for Tunisian Arabic", *LREC 2014*, p. 2355-2361.
- ZRIBI I., ELLOUZE M., HADRICH BELGUITH L., BLACHE P. (2015). "Spoken Tunisian Arabic Corpus "STAC": Transcription and Annotation", *Research in Computing Science 90*, p. 123-135.

Annexe 1 : convention de transcription pour les phénomènes oraux et interactionnels

Phénomène	Conventions	Exemples
Identité du participant		
Participant identifié	Toute parole est rapportée au locuteur qui l'a produite, qui est indiqué en début de ligne généralement par son identifiant, qui est le diminutif de son pseudonyme (utilisé pour des raisons éthiques et juridiques). L'identifiant est composé d'un, de deux ou de trois caractères (en majuscules). Dans les documents word, il est suivi d'une tabulation. Attention: c'est le seul endroit dans une transcription où l'usage de la tabulation est admis.	COR hāj
Doutes sur l'identité du participant	Lorsque le locuteur n'est pas identifiable, on utilise des X à la place du pseudo, en spécifiant F(emme) ou H(omme) si possible et nécessaire.	XXX hāj XXF bāj XXH salēm
Tour		
Notation du tour	La notation du tour est toujours rapportée à un identifiant. Un tour peut s'étendre sur plusieurs lignes. Le paragraphe d'un tour est aligné à gauche.	COR u baʕdēn hanballej duYri duYri #direct# waḷḷa
Enchaînement immédiat	= indique deux tours de parole qui s'enchaînent immédiatement sans micro-pause entre eux (<i>latching</i>).	COR kifek/= SAR =kifak anta
Chevauchement	Les crochets [et] encadrent le chevauchement dans chacun des tours concernés.	SAR hāj li[na] COR [ahlē:n]
Silences et pauses	Les pauses d'une durée inférieure à 0,2 secondes sont notées (.) (micro-pauses) Au-dessus, elles sont chronométrées à l'aide d'un logiciel au 10 ^{ième} de seconde près.	(.) (0.7)
Production de la parole		
Relation phonie-graphie	Pour les productions en arabe, on utilise une transcription phonético-phonologique (voir la convention ARAPI). Pour les productions dans d'autres langues, on utilise une transcription orthographique	SAM ahlēn fu xbārik SAR monsieur
Elision non standard	En cas de besoin, pour les langues transcrites en orthographe, on peut recourir à l'élision non standard, notée par une apostrophe ` (suivie d'un espace en fin de mot ; sans espace à l'intérieur du mot). Cette notation distinguant l'élision standard de l'élision non standard est importante pour le fonctionnement des outils de recherche automatique dans les données.	SAR c'est d'jà ça COR j'ai compris (.) j' recommence
Segments inaudibles	Les segments incompréhensibles sont représentés au moyen d'une série de x, chacun valant une syllabe. Si le nombre de syllabes ne peut être perçu, on note (inaud.) Si le transcritteur hésite entre plusieurs notations, les deux sont notées entre parenthèses, séparées par des points virgules.	COR xxx COR (inaud.) SAR (bʔūl; btəʔūl)
Allongement	Le son allongé est noté par : (sans espace avant). Les : sont répétés en fonction de la durée perçue de l'allongement.	COR wē:nek/ SAR hō:::ne ma feftini

Troncation	L'abandon d'un mot en cours de production est noté par - (sans espace avant, un espace après, ce qui le distingue du trait d'union).	COR ajm- ajmta beddek teʒi
Aspiration	L'aspiration est notée par .h (h précédé d'un point sans espace).	COR .h: eh:: lima la
Expiration, soupir	L'expiration est notée par h	COR h:: laʔa akid lā
Prosodie		
Montée et chute intonative	Les montées et chutes intonatives sont notées par « / » et « \ » sans espace avant. Les fortes montées et chutes sont notées par « // » et « \\ ».	SAM xallaʒte// LEI lā\ lissa
Saillance perceptuelle	Les segments perçus comme saillants sont soulignés.	SAM <u>waʒla</u> ma baʕref
Production vocale		
Production vocale	Les caractéristiques vocales et du mode de production sont décrites entre doubles parenthèses et précèdent la transcription, l'ensemble est compris entre chevrons :	COR <((en riant)) ma fi:ni>
Actions et événements. Cela concerne tous les phénomènes décrits et non transcrits, c'est-à-dire qui ne correspondent pas à des paroles produites par les participants.		
Action ou événement attribuable	Description entre doubles parenthèses dans le "pavé" du tour, après le pseudo du locuteur en début de ligne	CAR ((tape un numéro de téléphone))
Action ou événement	Même notation entre doubles parenthèses, mais sans identifiant dans la colonne des pseudos.	((le tonnerre retentit))

Annexe 2 : Liste des notations conventionnelles

0. PRESENTATION DES 4 TIERS.....	2
1. TRANSCRIPTION DES SONS.....	2
1. 1. CONSONNES.....	2
1. 2. VOYELLES BREVES.....	2
1. 3. VOYELLES LONGUES.....	3
1. 4. EXEMPLES.....	3
2. LONGUEUR PHONETIQUE ET EXPRESSIVE.....	4
3. VOYELLES LONGUES EN FIN DE MOT.....	4
4. TRANSCRIPTION DU ALIF SUSCRIT.....	4
5. TRANSCRIPTION DE LA HAMZA.....	4
6. LES DETERMINANTS.....	4
6. 1. L'ARTICLE.....	4
6. 2. LES DEMONSTRATIFS.....	4
7. LES PRONOMS.....	5
7. 1. LES PRONOMS PERSONNELS.....	5
7. 1. 1. <i>Les pronoms isolés sujets</i>	5
7. 1. 2. <i>Les pronoms suffixes possessifs</i>	5
7. 1. 3. <i>Les pronoms suffixes compléments</i>	5
7. 1. 4. <i>Exemples</i>	5
7. 2. LE PRONOM RELATIF.....	6
7. 3. LES PRONOMS INTERROGATIFS.....	6
7. 4. LES PRONOMS DEMONSTRATIFS.....	6
8. PARTICULES ET PREPOSITIONS.....	6
9. LES CHIFFRES.....	7
10. MORPHOLOGIE VERBALE.....	7
10.1. INACCOMPLI.....	7
10.2. INACCOMPLI MODAL.....	8
10.3. ACCOMPLI.....	8
10.4. IMPERATIF.....	8
10. 5. FORMES AUGMENTEES.....	8
10. 6. EXEMPLES.....	9
11. MORPHOLOGIE NOMINALE.....	9
12. AUTRES ELEMENTS CONVENTIONNELS.....	10
12. 1. LES EXPRESSIONS FIGEES (LISTE OUVERTE).....	10
12. 2. LES MOTS OU SYNTAGMES EN LANGUE ETRANGERE.....	10
12. 3. LES MOTS ETRANGERS ENTRES DANS LES USAGES (LISTE OUVERTE).....	10
12. 4. LES EMPRUNTS ARABISES (LISTE OUVERTE).....	11
12. 5. REGULATEURS, MARQUES D'HESITATION, ACQUIESCEMENT, NEGATION.....	11
12. 6. LES NOMS PROPRES.....	11
12. 7. LA MENTION DE NOMS ANONYMISES.....	11
13. LE LEXIQUE DIALECTAL (LISTE OUVERTE).....	11

0. Présentation des 4 tiers

La convention AraPI comprend 4 tiers :

Tier 1 : *transcription phonético-phonologique (API modifié)*

Tier 2a : *translittération morpho-phonologique (API modifié)*

Tier 2b : *glose morpho-syntaxique*

Tier 3 : *translittération arabe (avec caractères modifiés)*

Tier 4 : *traduction*

Dans la **tier 1**, on transcrit phonétiquement sans atteindre un degré de granularité très fin. Dans les tableaux qui suivent, les exemples donnés pour la tier 1 ne représentent pas toutes les réalisations possibles.

La **tier 3** est une notation conventionnelle du dialecte en caractères arabes : les traits morphologiques spécifiques au dialecte sont conservés, mais le *ductus* des mots se rapproche dès que possible de l'arabe littéral standard. De nombreux éléments courants, tels que les pronoms, les interrogatifs, etc. font l'objet de notations conventionnelles. La **tier 2a** est la translittération de la tier 3 en API.

La **tier 2b** est une glose (cf. annexe 4 pour une liste des gloses et des exemples).

La **tier 4** est une traduction dans la langue du chercheur, qui rend le sens de l'énoncé tout en s'attachant à rester au plus près de l'original.

1. Transcription des sons

1. 1. Consonnes

Les consonnes sont présentées dans l'ordre alphabétique de l'arabe.

Tier 3	Tier 1	Tier 2a	Tier 3	Tier 1	Tier 2a
ء	ʔ	ʔ	ط	t̤	t̤
ب	b, p	b	ظ	ð, d, z	ð
ت	t	t	ع	ɕ	ɕ
ث	θ, t, s	θ	غ	ɣ	ɣ
پ	p	p	ف	f, v	f
ج	ʒ, dʒ, g	ʒ	ڤ	v	v
ح	ħ	ħ	ق	q, g, ʔ	q
خ	x	x	ك	k, g	k
د	d	d	گ	g	g
ذ	ð, d, z	ð	ل	l, ɭ	l
ر	r	r	م	m	m
ز	z	z	ن	n	n
س	s	s	ه	h	h
ش	ʃ	ʃ	و	w	w
ص	ʂ	ʂ	ي	j	j
ض	d, z	z			

1. 2. Voyelles brèves

Il n'y a aucune vocalisation en tier 3. Pour les termes étrangers se terminant par le son [e], on peut le transcrire به en arabe.

Tier 1	Tier 2a
a, ə	a
ɛ, e, ə	e
i, ə	i
o, ə	o
u, y	u

1. 3. Voyelles longues

Tier 3	Tier 1	Tier 2a
ا	ā	ā
ا-ي	ē	ē
ي	ī	ī
و	ō	ū
و	ū	ū
ʔø	ʔø	ʔø (hésitation)
	ã	hã (marqueur exclamatif)
		bōʔūr

API	Exemple	API	Exemple
a	abjaɖ	ā	ktāb
ɛ	bərke	ē	lēʃ, bēt, bēred
e	ʔajjeb	ī	dīk
ə	sətt, təffēha	ō	ʃōb
i	kīfik	ū	sūʔ
o	bjektob	ã	ã... ok (exclamation)
u	suʔāl	ō	bōʔūr
y	styɖjo	ø	ʔø (hésitation)

1. 4. Exemples

Tier 1	Tier 2a	Tier 3
zyīr	ʃayīr	صغير
rfiʔa	rafīqa	رفيقة
tlēt	θalāθ	ثلاث
masalan	maθalan	مثلا
iza	iðā	إذا
bi ʔzabəʔ	bi-l-ɖabʔ	بالضبط
mazbūʔ	maɖbūʔ	مضبوط
zāhra	ðāhira	ظاهرة
ɖall	ðalla	ظل
aʎlā, aʎlāh	allāh	الله
bēt	bajt	بيت
jōm	jawm	يوم
ktēb	kitāb	كتاب
sijēse	sijāse	سياسة
baḥər	baḥr	بحر
azraʔ	azraq	أزرق
rās	raʔs	رأس
samā	samāʔ	سما

2. Longueur phonétique et expressive

Tier 1	Tier 2a	Tier 3
bʕi::d	baʕid	بعيد
lē:ʃ	lēʃ	ليس
aki::d	akid	أكيد

3. Voyelles longues en fin de mot

Tier 1	Tier 2a	Tier 3
ana	anā	أنا
ʕala	ʕalā	على
farʕine	farʕi-nī	فرجيني
fi	fī	في

4. Transcription du *alif* suscrit

Tier 1	Tier 2a	Tier 3
aḷḷā, aḷḷāh	allāh	الله
lākin, lakin	lākin	لكن
hāza, haza	hāḏā	هذا
halbēt	hāḏā l-bajt	هذا البيت

5. Transcription de la *hamza*

Tier 1	Tier 2a	Tier 3
bīr	biʔr	بئر
rās	raʔs	رأس
raʔis	raʔis	رئيس
samā	samāʔ	سما
tlēt ijjēm	ḡalāḡat ajjām	ثلاثة أيام
ḡkīne	iḡkī-nī	أحكي
esmi	ism-i	اسمي
ʃū əsmak	ʃū ism-ak	شو أسمك
xajje	ax-i	أخي

6. Les déterminants

6. 1. L'article

Tier 1	Tier 2a	Tier 3
à l'initiale		
elləʕbe	el-luʕba	العبة
eddār	el-dār	الدار
effaməs	el-ʃams	الشمس
elbejt	el-bajt	البيت
en liaison		
bə lləʕbe	bi-l-luʕba	بالعبة
bi ddār	bi-l-dār	بالدار
taḡt effaməs	taḡt l-ʃams	تحت الشمس
zēt ellōn	ḡāt l-lawn	ذات اللون
ʃakəl bēt elʕirān	ʃakl bajt l-ʕirān	شكل بيت الجيران
tannurt elbenet	tannurat l-bint	تنورة البنات

6. 2. Les démonstratifs

Tier 1	Tier 2a	Tier 3
--------	---------	--------

ha	ha	هـ
hajda, haza	hāḏā	هذا
hajde, hajdi, hazihi	hāḏī	هذي
hadōl, hadōle	haḏūl	هذول
hdāk	haḏāk	هذاك
hadīk, hadīke	haḏīk	هذيك

Tier 1	Tier 2a	Tier 3
halktēb bi halbēt	ha-l-kitāb bi-ha-l-bajt	مالكتاب بهالبيت

7. Les pronoms

7. 1. Les pronoms personnels

7. 1. 1. Les pronoms isolés sujets

Tier 1	Tier 2a	Tier 3
ana, ani	anā	أنا
enta	enta	إنت
ente	enti	إنتي
huwwe	huwa	هو
hijje	hija	هي
neḥna, ?ēḥna	naḥnā	نحنا
ento	entō	إنتوا
henne, hennen	hennen	هنن

7. 1. 2. Les pronoms suffixes possessifs

Tier 1	Tier 2a	Tier 3
ktēbe	kitāb- i	كتابي
ktēbek / ktēbak	kitāb- ka	كتابك
ktēbek / ktēbik	kitāb- ki	كتابكي
ktēbo	kitāb- hu	كتابه
ktēba	kitāb- hā	كتابها
ktēbna	kitāb- nā	كتابنا
ktēbkon	kitāb- kon	كتابكن
ktēbon	kitāb- hon	كتابهن

7. 1. 3. Les pronoms suffixes compléments

Tier 1	Tier 2a	Tier 3
jāni, jāne	jē- nī	إياني
jāk, ijjāk, jēk	jē- ka	إياك
jāki, ijjāki	jē- ki	إياكي
jē, ijjē, jāho	jē- hu	إياه
jē, ijjē, jāha	jē- hā	إياها
jāna	jē- nā	إيانا
jekon	jē- kon	إياكن
jon / jāhon	jē- hon	إياهن

7. 1. 4. Exemples

Tier 1	Tier 2a	Tier 3
bētna, bajtna katbētha, katabēta ʕaṭīne, ʕṭīne ?āletle	bajt-nā katab-at-hā aʕṭi-nī, aʕṭī-nī qāl-at l-i	بيتنا كتبتها اعطني / اعطيني قالت لي
biddi, bēddi bēddak, biddak, baddak bēddik, biddik, baddik	badd-i badd-ka badd-ki	بدي بذك بذكي

bəddo	badd-hu	بده
bədda	badd-hā	بدها
badna	badd-nā	بدهنا
bədkon	badd-kon	بدهكن
bəddon	badd-hon	بدهن

7. 2. Le pronom relatif

Tier 1	Tier 2a	Tier 3
jilli, elli, li, l	jalli	اللي

7. 3. Les pronoms interrogatifs

Tier 1	Tier 2a	Tier 3
kīf	kīf	كيف
wejn, wēn	wīn	وين
fejn, fēn	fīn	فين
mīn	mīn	مين
lē	lē	ليه
lēʃ	lēʃ	ليش
ēʃ	ēʃ	ايش
ʃū	ʃū	شو
ʔaddēʃ	qaddēʃ	قديش
ʔadde	qaddē	قديه
ejmata	ejmatā	ايمتى
emta	emtā	ايمتى
ej / ejja	aj / aja	أي / أية

7. 4. Les pronoms démonstratifs

Tier 1	Tier 2a	Tier 3
hajda, haza	hāḏā	هذا
hajde, hajdi, hazihi	hāḏī	هذي
hadōl, hadōle	haḏūl	هذول
hdāk	haḏāk	هذالك
hadīk, hadīke	haḏīk	هذيك

8. Particules et prépositions

Toutes les particules et prépositions sont séparées par une espace en tier 1. En tier 2a, les particules monolitères sont segmentées par un tiret. En tier 3, on respecte les règles d'orthographe de l'arabe standard.

Tier 1	Tier 2a	Tier 3
Particules monolitères		
bi ʃʃēf	bi-l-ʃajf	بالصيف
bi lbēt	bi-l-bajt	بالبيت
bi bēt elʒirān	bi-bajt l-ʒirān	ببيت الجيران
bi lxzēne	bi-l-xizāne	بالخزانة
ta jektob	ta-ja-ktub	ت يكتب
la jkūn	li-ja-kūn	ليكون
la lmadām	li-l-madām	للمدام
ʕa tṭawle	ʕa-l-ṭāwile	عاطولة
Autres particules		
mn əlbīdāje	min l-bīdāje	من البداية
foʔ ərrās	fawq l-raʔs	فوق الرأس

9. Les chiffres

Les chiffres sont transcrits phonétiquement en tier 1, mais notés conformément à leur équivalent en arabe standard dans les tiers 2a et 3.

Tier 1	Tier 2a	Tier 3
ʃefer	ʃifr	صفر
wāḥad	wāḥid / aḥad / iḥdā	واحد / إحدى
tnēn	iḥnajn / iḥnatajn	إثنين / إثنين
tlēte	ḥalāḥ / ḥalāḥa	ثلاث / ثلاثة
arbʃa	arbaʃ / arbaʃa	أربع / أربعة
xamse	xams / xamsa	خمس / خمسة
sette	sitt / sitta	ست / ستة
sabʃa	sabʃ / sabʃa	سبع / سبعة
tmēne	ḥamānin / ḥamāniya	ثمان / ثمانية
tesʃa	tisʃ / tisʃa	تسع / تسعة
ʃaʃra	ʃaʃar / ʃaʃra	عشر / عشرة
ḥdaʃʃ	aḥad ʃaʃar / iḥdā ʃaʃra	أحد عشر / إحدى عشرة
tnaʃʃ	iḥnaj ʃaʃar / iḥnataj ʃaʃra	إثني عشر / إثنى عشر
tlattaʃʃ	ḥalāḥat ʃaʃar	ثلاثة عشر
ʃarbaʃtaʃʃ	arbaʃat ʃaʃar	أربعة عشر
xamstaʃʃ	xamsat ʃaʃar	خمس عشر
settaʃʃ	sittat ʃaʃar	ستة عشر
sabaʃtaʃʃ	sabʃat ʃaʃar	سبعة عشر
tmantaʃʃ	ḥamānijat ʃaʃar	ثمانية عشر
tesaʃtaʃʃ	tisʃat ʃaʃar	تسعة عشر
ʃeʃrīn	ʃiʃrīn	عشرين
tlētīn	ḥalāḥīn	ثلاثين
arbʃīn	arbaʃīn	أربعين
xamsīn	xamsīn	خمسین
settīn	sittīn	ستين
sabʃīn	sabʃīn	سبعين
tmēnīn	ḥamānīn	ثمانين
tesʃīn	tisʃīn	تسعين
mit	miʔa	مئة
eləf	alf	ألف

10. Morphologie verbale

Les désinences verbales ne sont segmentées par un tiret que dans le tier 2 où le radical du verbe est isolé de l'ensemble des affixes (morphèmes de modes/temps et personnes) qui font l'objet d'une notation conventionnelle.

Attention : la convention s'écarte de la proposition de CODA dans la notation de la conjugaison à la 1^{ère} personne du pluriel (on note m-n- et pas b-n-).

En cas d'élimination du morphème de conjugaison en tier 1 ('am fakkir), on rétablit en tier 2 et 3 le morphème manquant ('am u-fakkir). La voyelle du radical et des personnes (pour les formes dérivées) est transcrite phonétiquement en tier 1 mais devient conventionnelle en tier 2.

Voici les tableaux de conjugaison pour le verbe *katab*, suivis de différents exemples pour d'autres verbes.

10.1. Inaccompli

Tier 1	Tier 2a	Tier 3
bektob	b-a-ktub	بأكتب
btektob	b-ta-ktub	بنكتب
btektobe	b-ta-ktub-i	بنكتبني
bjektob	b-ja-ktub	بيكتب

btektob	b- ta -ktub	بنتكتب
mnektob	m- na -ktub	منكتب
btektobo	b- ta -ktub-ū	بنتكتبوا
bjektobo	b- ja -ktub-ū	بيكتبوا

10.2. Inaccompli modal

Pour l'inaccompli à valeur modale (sans le préverbe b-) employé après différentes particules comme حتى, لازم, بدّ, في, etc., la hamza de la première personne du singulier est rétablie dans les tiers 2 et 3.

Tier 1	Tier 2a	Tier 3
ektob	a -ktub	أكتب
tektob	ta -ktub	تكتب
tektobe	ta -ktub-ī	تكتبي
jektob	ja -ktub	يكتب
tektob	ta -ktub	تكتب
nektob	na -ktub	نكتب
tektobo	ta -ktub-ū	تكتبوا
jektobo	ja -ktub-ū	يكتبوا

10.3. Accompli

Tier 1	Tier 2a	Tier 3
katabet, katabt	katab- t	كتبت
katabet, katabt	katab- t	كتبت
katabte	katab- tī	كتبتني
katab	katab	كتب
katbet	katab- at	كتبت
katabna	katab- nā	كتبنا
katabto	katab- tū	كتبوا
katabo	katab- ū	كتبوا

10.4. Impératif

Tier 1	Tier 2a	Tier 3
ktōb	u-ktub	اكتب
ktebe	u-ktub-ī	اكتبي
ktebo	u-ktub-ū	اكتبوا

10.5. Formes augmentées

	Tier 1	Tier 2a	Tier 3
II	ħammam	ħammam, ju-ħammim	حَمَمَ ، يَحَمَمُ
	ʔattal	qattal, ju-qattil	قَتَلَ ، يَقْتُلُ
III	ʕāmal	ʕāmal, ju-ʕāmil	عَامَلَ ، يَعْمَلُ
	kētab	kētab, ju-kātib	كَاتَبَ ، يَكْتُبُ

IV	aḏhar	aḏhar, ju-ḏhir	أظهر ، يظهر
V	tḥammam	tḥammam, ja-ta-ḥammam	تَهَمَّ ، يتَهَمَّ
	tkassar	tkassar, ja-ta-kassar	كَسَّرَ ، يَتَكَسَّرُ
VI	tṣāmal	tṣāmal, ja-ta-ṣāmal	تَعَامَلَ ، يَتَعَامَلُ
	tkētaf	tkētaf, ja-ta-kātaf	كَاتَفَ ، يَتَكَاتَفُ
VII	nkasar	nkasar, ja-n-kasir	انكسر ، ينكسر
	nʔaṭaṣ	nqataṣ, ja-n-qatiṣ	نَقَطَعَ ، يَنْقَطِعُ
VIII	ʃtaḃal	ʃtaḃal, ja-ʃtaḃil	اشتغل ، يشتغل
	ṣtamad	ṣtamad, ja-ṣtamid	اعتمد ، يعتمد
X	staʔzan	staʔzan, ja-sta-ʔzin	ستازن ، يستازن
	stayfar	stayfar, ja-sta-ʔfir	ستغفر ، يستغفر

10. 6. Exemples

Tier 1	Tier 2a	Tier 3
minʔassemon	m-nu-qassim-on	منقسمهن
jkūn / jikūn	ja-kūn	يكون
ṣam jektob	ṣam ja-ktub	عم يكتب
raḥ tektob	raḥ ta-ktub	رح تكتب
ḥa jektob	ḥa-ja-ktub	حي يكتب
ṣam fakkər	ṣam ufakkir	عم أفكر
bteṣrabe	b-ta-ṣrab-ī	بتشربى
jaṣmol	ja-ṣmal	يعمل
semṣet	samiṣ-t	سمعت
balleṣ	ballaṣ	بلش
biballeṣ	b-ju-balliṣ	بيبلش
betḥebb	b-tu-ḥibb	بتحب
mən rūḥ	m-na-rūḥ	منروح
fallet	fall-at	فالت
nsīna	nasī-nā	نسينا
ḥajaṣtūne	ḥa-ja-ṣt-ū-nī	حيعطوني
xallīni nām	xallī-nī anām	خليني أنام
ʔāletlon	qāl-at la-hon	قالت لهن
ʔaltella	qul-t la-hā	قالت لها
biddi ḥʔuz	badd-ī aḥʔuz	بدي أحجز

11. Morphologie nominale

Dans la tier 1, les marques de genre et nombre des participes et des pluriels externes ne sont pas segmentées. Dans la tier 2a, elles sont segmentées. Conventionnellement, le pluriel externe masculin est noté -in.

Tier 1	Tier 2a	Tier 3
ṣāmil	ṣāmil	عامل
ṣāmle	ṣāmil-e	عاملة
ṣāmlīn	ṣāmil-in	عاملين
mṣallimīn	muṣallim-in	معلمين
mṣallimāt	muṣallim-āt	معلمات

12. Autres éléments conventionnels

12. 1. Les expressions figées (liste ouverte)

Voici la liste des expressions qui sont codées en tier 2 et 4 et qui font l'objet d'une notation conventionnelle en tier 3, par ordre alphabétique arabe.

Tier 1	Tier 2a et 4	Tier 3
akīd	AKID	أكيد
in fā ʔaḷḷāh, nfaḷḷa	INCH	إن شاء الله
enno	ENNO	إنّه
bass, bæss	BASS	بس
hallāʔ	HALQ	هلاق
lḥamdəḷḷa	HDLA	الحمد لله
aḷḷa ʔirḥamo	ARHM	الله يرحمك
lakēn	LKAN	لكان
lajke	LAYK	ليك
mafāḷḷa	MCHA	ما شاء الله
waḷḷā	WALA	والله
u hēke	WHEK	وهاك
jaḷḷā	YALA	يالله
jaʕne, jaʕni	YANI	يعني

Tier 1	Tier 2a	Tier 3	Tier 4
lḥamdəḷḷa mēfe lḥāl kifīk ʔnte	HDLA māfī l-ḥāl kif-kī enti	الحمد لله ماشي الحال كفيك إنتي	HDLA ça va et toi
lajka:/ kənnā badna hēk naʕmol #réunion#	LAYK kun-nā badd- nā hāk na-ʕmal #réunion#	ليك كنا بدنا هاك نعمل #réunion#	LAYK on voulait comme ça faire une #réunion ((FR))#

12. 2. Les mots ou syntagmes en langue étrangère

Tier 1	Tier 2a	Tier 3	Tier 4
bōʕūr ṣabāja sori ajja #facture# la ljōm #dernier délai#	#bonjour# ṣabāja #sorry# ajja #facture# li- l-jawm #dernier délai#	#بنجور# صبايا #سوري# أية #facture# لليوم #dernier délai#	#bonjour ((FR))# les filles #désolée ((GB))# quelle #facture ((FR))# est pour aujourd'hui #dernier délai ((FR))#
enno Nicolas #hair# baddo #bonnet#	ENNO Nicolas #hair# badd-hu #bonnet#	إنّه نقولا #hair# بده #bonnet#	ENNO Nicolas a besoin d'un #bonnet ((FR))# pour ses #cheveux ((GB))#

12. 3. Les mots étrangers entrés dans les usages (liste ouverte)

Voici des exemples de mots d'origine étrangère entrés dans les usages : ils figurent en tier 2a dans l'orthographe de la langue. C'est une liste ouverte qui s'enrichira au fur et à mesure des transcriptions de nouveaux corpus.

Tier 1	Tier 2a	Tier 3	Tier 4
bonʕūr / bōʕūr sori / soyi mersi sa va hāj	#bonjour# #sorry# #merci# #ça va # #hi#	#بنجور# #سوري# #مرسي# #ساقا# #هاي#	#bonjour ((FR))# #désolé (e) ((GB))# #merci ((FR))# #ça va ((FR))# #salut ((GB))#

bāj oke	#bye# #okay#	#باي# #أوكي#	#au revoir ((GB))# #okay ((GB))#
------------	-----------------	-----------------	-------------------------------------

12. 4. Les emprunts arabisés (liste ouverte)

Voici des exemples d'emprunts arabisés : ils figurent en tier 1 et 2a en transcription phonétique. C'est une liste ouverte qui sera enrichie au fur et à mesure des transcriptions de nouveaux corpus.

Tier 1	Tier 2a	Tier 3	Tier 4
#pojēt# #taperwarēt#	#pojāt# #taperwarāt#	#پويات# #تپيروارات#	#des pots ((FR))# #des tupperwares ((GB))#

12. 5. Régulateurs, marques d'hésitation, acquiescement, négation

Sens	Tier 1	Tier 2a	Tier 3	Tier 4
clic = non	tsk	tsk	tsk	tsk
morphème de négation	həʔə	həʔə	həʔə	həʔə
oui	ʔē	ʔē	إيه	oui
accusé réception minimal	ʔm	ʔm	ʔm	ʔm
accusé réception double	mhm	mhm	mhm	mhm
tag	ha	ha	ha	ha
morphème d'exclamation	a::::h, ā	āh	أه	ah
hésitation	ʔø	ʔø	ʔø	euh
hésitation (mot tronqué)	k- ktīr	k- kaθīr	ك كثير	beau- beaucoup

12. 6. Les noms propres

Les noms propres (anonymisés ou non) figurent entre guillemets hauts dans les trois premières tiers. Il peut s'agir d'une personne, d'un lieu, d'une institution, etc.

Tier 1	Tier 2a	Tier 3	Tier 4
mitəl "Woody Allen"	miθl "Woody Allen"	مثل "وودي آلن"	comme Woody Allen
aflām "Jusef jahin" elʔadīme	aflām "Jūsuf fāhin" el-qadīma	أفلام "يوسف شاهين" القديمة	les vieux films de Youssef Chahine

12. 7. La mention de noms anonymisés

Dans toutes les tiers de la transcription, les noms propres anonymisés cités apparaissent sous leur forme intégrale ou abrégée selon les cas. Cela peut concerner une personne, un lieu, une institution, etc. S'il s'agit du pseudo de l'un des participants à l'interaction, on le conserve sans guillemets.

Tier 1	Tier 2a	Tier 3	Tier 4
marħaba SOF (0.5) kif-ik	marħaba SOF (.) kif-ik	مرحبا SOF (.) كيفك	Salut SOF (.) ça va
marħaba ja "3or3"	marħaba jā "Georges"	مرحبا يا "جورج"	Salut Georges

13. Le lexique dialectal (liste ouverte)

Dans ce tableau, amené à être enrichi, on a relevé les mots dialectaux et expressions très courants ou sans équivalents en arabe standard. La tier 1 présente une liste des variantes alors que les tiers

2a et 3 proposent la notation conventionnelle correspondante. Les termes figurent dans l'ordre alphabétique arabe de la tier 3.

Tier 1	Tier 2a	Tier 3	Tier 4
alo, a:lo	alo	ألو	Allô
oke, okaj	oke	#أوكي#	okay
ənno	ʔnno	انه	que
jā... jā	jā... jā	أو	ou
wala	wa-lā	أو	ou
ūḍa	ūḍā	أوضى	chambre, pièce
awwal ma	awwal mā	اول ما	dès que
ajwa	ajwā	أيوا	oui, bien sûr
ḅāj	bāj	#باي#	salutation de clôture
beḥkiki, mneḥke	ba-ḥkī-ik, m-na-ḥkī	بحكيكي, منحكي	je te parle, on se parle (salutation de clôture, à bientôt)
bfūfek nfālla	ba-fūf-ik INCH	بشوفك إن شاء الله	je te vois INCH (salutation de clôture, à bientôt)
baddak / biddak badda	badd-ka badd-hā	بدك بدها	tu veux elle veut
beddi ʔəzəʔʔek	badd-ī u- zʔiʔ-ak	بدي از عك	je vais te/vous déranger
baddik fī	badd-ik fī	بدك شي	tu veux quelque chose ? (pré-salutation de clôture)
barra	barra	برّا	à l'extérieur
barrat...	barrat...	برّاة...	à l'extérieur de...
barḍo	barḍo	برضو	encore, aussi
bass	bass	بس	quand, seulement, juste
b ʔaleb	bi-qalb	بقلب	dans, à l'intérieur
bala	bilā	بلا	sans
balke, barke	balke	بلكي	peut-être
tabaʔ tabʔet tabʔūl	tabaʔ tabʔet tabʔūl	تبع تبعه تبعول	à, de
ʔuwwa	ʔuwwa	جوا	à l'intérieur, dedans
ʔuwwat...	ʔuwwat...	جواة...	à l'intérieur de...
ḥadd	ḥadd	حد	à côté de
lḥamdəlla: mefe lḥāl	HDLA māfī l-ḥāl	الحمد لله ماشي الحال	HDLA ça va
lḥamdəlla tamēm	HDLA tamām	الحمد لله تمام	HDLA très bien, super
dəyri	duyri	دغري	directement
rāḥ, raḥ	raḥ	رح	(particule du futur)
zalame	zalame	زلمة	gars
zejj	zajj	زي	comme
salām	salām	سلام	(salutation d'ouverture)
ʃwaj, ʃwajj, ʃwej	ʃwaj	شوي	un peu
ʃwajet...	ʃwajat-	شوية...	un peu de...
ʃī	ʃī	شي	(quelque) chose
ṭajjeb	ṭajjib	طيب	bon, bien

ʕafān, ʕalafān	ʕafān	عشان	pour
fəm, təm	fam	فم	bouche
karmēl	karmāl	كرمال	pour
kamēn, kamēna	kamān	كمان	aussi
lā	lā	لا	non
laʔ, laʔa	lāʔa	لأ	certainement pas
lessa, essa	lissa	لسة	pas encore
mbārēh	mbārih	مبارح	
mbala, əmbala	mbalā	مبلا	non, au contraire
madām, madāmtak	madām, madāmat-ka	مدام	madame
madri	madrī	مدري	je ne sais pas
məʃ, miʃ, muʃ	mʃ	مش	(négation)
meʃtēʔetlek	muʃtāqa la-ki	مشناقة لك	tu me manques
miʃān, minʃān	miʃān	مشان	pour
mēʃe lħāl	māʃi l-ħāl	ماشي الحال	ça va
mnīh, mliħ	mnīh	منيح	bien
mū	mū	مو	(négation)
nijjāl...	nijjāl-	نيال	tant mieux pour...
hallaʔ	hallaq	هلاق	maintenant
hōn, hōne	hunā	هنا	là
honik, hnēke	hunāk	هناك	là-bas
hēk	hāk	هيك	comme ça
hāj	hāj	هاي	hi (salutation d'ouverture)
jalla	jallāh	يللي	allez
jalla ʔāj	jallāh bāj	يللي باي	allez salut

Annexe 3 : Liste indicative des gloses morpho-syntaxiques

1	première personne	IMP	impératif
2	deuxième personne	INDF	indéfini
3	troisième personne	INTENS	intensifieur
ACC	accusatif	INTER	interrogatif
ACP	accompli	LOC	locatif
ACT	actif (vs. passif)	M	masculin
ADJ	adjectif	N	nom
ADV	adverb(ial)	NEG	négation
AFX	affixe	NPR	nom propre
AGR	accord, agreement	NOM	nominatif
ANN	annexion	PART	particule
ART	article	PAST	passé
AUX	auxiliaire	PRE	préposition
CHI	nombre cardinal ou ordinal	PAS	passif
COLL	collectif	PFX	préfixe
COMP	complémentiseur (ex. ENNO, que complétif)	PL	pluriel
COP	copule (pronominal copula)	POSS	possessif
DAT	datif	PRO	pronom personnel
DEF	défini	PRV	préverbe
DEM	démonstratif	PRES	présent
DIM	diminutif	PTA	participe actif
DU	duel	PTP	participe passif
ELAT	élatif	REL	relatif
F	féminin	SFX	suffixe
FUT	futur	SG	singulier
INA	inaccompli	SUPER	superlatif
IFX	infixe	V	verbe
		VOC	vocatif

Annexe 4 : Transcrire les dialectes arabes en linguistique. Comparaison des conventions existantes

Document réalisé par Catherine PINON

La transcription des dialectes arabes en caractères arabes est problématique, car il n'existe aucune orthographe officielle et les différences phonétiques et morpho-phonologiques peuvent être très importantes d'un dialecte à l'autre. Proposer une convention de transcription unique adaptable à tous les dialectes faciliterait les recherches et leur diffusion. Plusieurs propositions ont été faites dans ce sens. Nous avons choisi ici d'en présenter trois, sur lesquelles reposent les choix de transcription effectués dans la convention AraPI : les conventions CODA, AMADAT et BBN. Nous mentionnerons les règles propres à AraPI quand elles s'écartent des propositions présentées ici.

La réflexion qui suit sur la transcription des dialectes en caractères arabes ne vise pas l'établissement d'une orthographe standardisée. Elle propose une comparaison de systèmes existants, analysant ce qui les rapproche et ce qui les différencie. Il faut garder en tête que ces conventions de transcription ont pour but de proposer une référence bien exhaustive et bien explicitée. Son objectif majeur est de permettre le partage des données entre des chercheurs de différentes disciplines.

Ce document vise, dans le même esprit que les *Leipzig Glossing Rules*, à déduire des différentes propositions faites des règles communes, afin d'uniformiser les pratiques de transcription des dialectes. Après l'énoncé de règles générales, on prendra l'exemple du dialecte syro-libanais pour illustrer la méthode de transcription à suivre.

Table des matières

1. LES TRAVAUX EXISTANTS	2
1. 1. LA CONVENTION CODA	2
1. 2. LA TRANSCRIPTION AMADAT	2
1. 3. LES DIRECTIVES BBN	2
2. LES REGLES DE TRANSCRIPTION COMMUNES	3
REGLE 1 : UN MOT = UNE SEULE ET UNIQUE TRANSCRIPTION.....	3
REGLE 2 : UNE ORTHOGRAPHE AYANT COMME REFERENCE LE MSA.....	3
REGLE 3 : UNE MORPHO-SYNTAXE QUI PRESERVE LES SPECIFICITES DIALECTALES.....	4
REGLE 4 : ADAPTATION, MAIS COHERENCE ET PERTINENCE	5
3. LES REGLES DE TRANSCRIPTION SPECIFIQUES.....	5
REGLE 5 : LA TRANSCRIPTION DES MOTS ETRANGERS.....	5
REGLE 6 : CODE-SWITCHING ET DIGLOSSIC CODE-SWITCHING.....	6
REGLE 7 : NOMS PROPRES ET ABBREVIATIONS	6
4. QUESTIONS COURANTES.....	6
5. BIBLIOGRAPHIE COMPLEMENTAIRE.....	8

1. Les travaux existants

Les linguistes travaillant sur l'arabe sont tôt ou tard confrontés à la question de la transcription de leurs sources. Chaque chercheur fait ses propres choix en fonction de ses objectifs scientifiques, mais aussi de contraintes techniques qui peuvent peser sur ses recherches. Il existe plusieurs conventions qui visent à unifier les systèmes de transcription utilisés par les linguistes arabisants.

1. 1. La convention CODA

La convention CODA (Conventional Orthography for Dialectal Arabic) a été établie par Nizar Habash, Mona Diab et Owen Rambow du Center for Computational Learning Systems de l'Université Columbia à New York, en 2011.

> *références en ligne :*

Habash, Diab et Rambow, "Unified Guidelines and Resources for Arabic Dialect Orthography" [En ligne] <http://www.lrec-conf.org/proceedings/lrec2018/pdf/395.pdf>

Habash, Diab et Rambow, "Conventional Orthography for Dialectal Arabic" [En ligne] http://www.lrec-conf.org/proceedings/lrec2012/pdf/579_Paper.pdf

1. 2. La transcription AMADAT

Il s'agit d'une transcription dialectale à base orthographique (Dialectal Arabic Orthography-based transcription) mise au point par Mohamed Maamouri, David Graff, Hubert Jin, Christopher Cieri et Tim Buckwalter du Linguistic Data Consortium de l'Université de Pennsylvanie, pour l'étude des conversations téléphoniques (Conversational telephone speech = CTS) en arabe levantin.

> *références en ligne :* Maamouri, Graff, Jin, Cieri et Buckwalter, "Dialectal Arabic Orthography-based Transcription & CTS Levantine Arabic Collection" [En ligne] <https://www ldc.upenn.edu/sites/www ldc.upenn.edu/files/ears-rt04-dialectal-arabic-transcription.pdf>

1. 3. Les directives BBN

Les directives pour transcrire les dialectes arabes (Guidelines for Transcribing Arabic Dialects) ont été proposées par Bushra Zawaydeh, Dave Stallard et John Makhoul de BBN Technologies en 2002-2003.

> *références en ligne :* Zawaydeh, Stallard et Makhoul : "Guidelines for Transcribing Arabic Dialects" [En ligne] <https://catalog ldc.upenn.edu/docs/LDC2005S08/BBN-Babylon-transcription-guidelines.pdf>

2. Les règles de transcription communes

D'une manière générale, les transcriptions se veulent à la fois panarabiques et spécifiques. Elles visent la facilité d'emploi et la lisibilité. À partir d'un examen détaillé de ces conventions, nous avons distingué les éléments communs et les règles spécifiques qui s'en dégagent. Nous en proposons un résumé sous la forme d'un ensemble de règles.

Règle 1 : un mot = une seule et unique transcription

Chaque mot a une unique forme orthographique en transcription, même si sa réalisation phonétique varie.

Par conséquent, on peut établir pour chaque mot une transcription conventionnelle et une liste des variantes.

Règle 2 : une orthographe ayant comme référence le MSA

Dès que cela est possible, les choix de transcription sont conformes à l'orthographe de l'arabe littéral (Modern Standard Arabic ou MSA), qui est assez stable sur l'ensemble du monde arabe et constitue une référence partagée par tous les arabophones scolarisés. Par exemple, la cliticisation des particules monolitères, le recours à la *šadda* pour la gémiation phonologique, l'écriture de l'article ou l'emploi du *tā' marbūṭa* et du *alif maqṣūra* sont conformes à l'orthographe standard.

Par conséquent, on cherche à établir pour chaque mot un lien avec le MSA.

- Si le mot a une forme directement apparentée en MSA, utiliser l'orthographe MSA (ce qui peut donner une forme qui n'existe pas en MSA, cf. exemple 4).

- Exemples :
- | | |
|--|--|
| (1) Libanais
/kti:R/ ou /kati:R/
كثير (non pas كتير) | <i>l'interdentale [θ] réalisée [t] est rétablie</i> |
| (2) Libanais
/daʔən/
دفن (non pas دفن ou دان) | <i>le [ð] et le [q] réalisés [d] et [ʔ] sont rétablis</i> |
| (3) Libanais
/zɣīr/
صغير (non pas زغير) | <i>le [ɣ] réalisé [z] est rétabli</i> |
| (4) Libanais
/hadōl/
هدول (non pas هدول) | <i>le [ð] des démonstratifs standard est rétabli, même si cette forme n'existe pas en MSA.</i> |

- Dans le cas où la racine commune au dialecte et au MSA a une consonne ajoutée ou élidée, on rétablit aussi la racine MSA.

Exemples : (5) Libanais
/u-noʃ/
ونصف (non pas نص)

- On réintègre les voyelles longues de schèmes particuliers même si elles ne sont pas réalisées et on note les tanwīn-an¹.

Exemples : (6) Libanais
/ṭabūr/
طبور (non pas طور) le [ā] réalisé [a] est rétabli

(7) Libanais
/matalan/
مثلاً (non pas مثلن) le [θ] réalisé [t] est rétabli et [an] est noté [ʾ]

- Si on utilise des textes déjà transcrits en dialecte, on corrige les erreurs typographiques (métathèses, omission ou ajout d'espaces). Par exemple, on rétablit la bonne orthographe du yā' / alif maqṣūra dans les textes égyptiens.

- Dans le cas où le lien entre la forme dialectale et le MSA n'existe pas, on livre une liste des transcriptions conventionnelles.

Exemples : (8) Libanais
/doʁri/
دغري (notation conventionnelle)

Règle 3 : une morpho-syntaxe qui préserve les spécificités dialectales

Bien que s'inspirant pour la transcription des sons du *ductus* classique, les spécificités morpho-syntaxiques des dialectes sont maintenues.

Par conséquent, on conserve les marques de temps et d'aspect dans la conjugaison, les terminaisons des pluriels externes, les pré-verbes dialectaux (par exemple, on ne remplace pas le marqueur de futur ḥa- par sa-) et on conserve les spécificités lexicales dialectales.

Exemple : (9) Libanais
/b-ya-ktub/
بيكتب (non pas كتب)

(10) Tunisien
/bdī-t/
بدیت (non pas بدأت)

(11) Libanais
/raħ balliʃ/
رأح أبأش (non pas سأبدأ)

(12) Libanais
/fêto ləmʕalmīn/
دأحل المألمون (non pas فأأأوا المألمین)

Règle 4 : adaptation, mais cohérence et pertinence

Les chercheurs adaptent les conventions existantes en fonction des objectifs de leur recherche. Ils doivent alors expliciter leurs choix. Dans tous les cas, les conventions établies doivent être suivies strictement et les exceptions listées.

3. Les règles de transcription spécifiques

En fonction des objectifs de la recherche, certains choix doivent être faits concernant des notations particulières.

Règle 5 : la transcription des mots étrangers

Les systèmes CODA, AMADAT et BBN s'accordent sur le fait de transcrire les mots étrangers en écriture arabe sans recourir à des caractères supplémentaires comme چ, پ ou ف. On utilise la transcription arabe usuelle si elle existe. Dans le cas contraire, on peut modifier la valeur de certaines lettres pour ne pas recourir à des caractères supplémentaires :

- on utilise le *bā'* (ب) pour transcrire [p]
- on utilise le *tā'+šīn* (تش) pour transcrire [tch]
- on utilise le *fā'* (ف) pour transcrire [v]
- on utilise le *jīm* (ج) ou le *kāf* (ك) pour transcrire [g]

Règle 5a (optionnelle) : préséance des formes levantines

En cas de concurrence entre plusieurs transcriptions connues, CODA préfère à la transcription égyptienne la transcription levantine : on écrit (كراج), que l'on transcrive de l'égyptien ou du libanais.

Règle 5b (optionnelle) : adaptation au standard dialectal

En cas de concurrence entre plusieurs transcriptions connues, AMADAT propose de se conformer avec les usages locaux : on écrit (كراج) si l'on transcrit du libanais et (جراج) si l'on transcrit de l'égyptien.

Règle 5c (optionnelle) : caractères arabes supplémentaires

AraPI recourt, dans le cas des emprunts arabisés, à un certain nombre de caractères supplémentaires, tels que ك, پ ou ف, en indiquant à l'aide de l'alphabet phonétique international leur valeur.

Règle 5d (optionnelle) : transcription en caractères latins

AraPI propose de transcrire les segments en langue étrangère en caractères latins, dans leur orthographe d'origine.

Règle 6 : code-switching et diglossic code-switching

Le marquage du passage d'une langue à une autre n'apparaît pas toujours dans les différentes conventions.

Règle 6a (optionnelle) : uniformisation (dialectalisation du MSA)

Dans CODA, si un terme est prononcé en MSA, on le transcrit dans le dialecte. Par exemple, si le mot « homme » est prononcé [raʒuɫ] comme en MSA dans un discours en égyptien, il sera transcrit راجل [rāgəɫ] conformément à la spécificité lexicale de ce mot en égyptien.

Règle 6b (optionnelle) : marquage du code-switching

Dans AraPI, les mots ou syntagmes en langue étrangère ou dans une variété d'arabe différente figurent encadrés de #.

Règle 7 : noms propres et abréviations

La transcription des noms propres ou des abréviations peuvent poser problèmes. Dans l'ensemble des conventions, on suggère d'utiliser la transcription consensuelle quand elle existe. D'autres propositions sont parfois faites concernant spécifiquement les noms propres.

Règle 7a (optionnelle) : marquage des noms propres et abréviations

BBN propose de lier par un tiret bas les noms propres et abréviations composées, par exemple :

برج_البراجنة

Règle 7b (optionnelle) : marquage des noms propres et abréviations

AraPI propose d'entourer la transcription des noms propres par des guillemets. Exemple :
"ماركس"

4. Questions courantes

Les éléments linguistiques les plus courants ne font pas toujours l'objet de règles détaillées. Voici point par point les propositions partagées par les conventions (quand elles les abordent), avec la mention des divergences quand le consensus n'est pas établi.

Les pronoms isolés et suffixes

Chaque dialecte a son écriture conventionnelle des pronoms. La 2^{ème} personne du féminin singulier est toujours marquée par un *yā'* : إنتي - كي

Pour les pronoms commençant par un *hā'*, on le note toujours même s'il n'est pas prononcé.

Les verbes

- on garde l'allongement vocalique des redoublés et des hamzés de 3^{ème} radicale (بديت)
- on note la 2^{ème} personne du féminin singulier par un *yā'* (بديتي)
- on note la 3^{ème} personne du masculin pluriel par وا (le *alif* n'est pas noté uniquement s'il y a un suffixe après : بيكتبوا - بيكتبوها)

Les préverbes sont attachés au verbe s'ils sont monolitères, sinon ils sont séparés par une espace, conformément à l'écriture arabe.

حأكتب - تيروح - عم بنكتب - رح بيكتبوا

La variante négative شي du négatif ش est toujours orthographiée ش à la fin du verbe.

Variantes : Pour CODA, le paradigme de conjugaison du libanais est :

بأكتب - بنكتب - بتكتبي - بيكتب - بتكتب - بنكتب - بتكتبوا - بيكتبوا

Pour BBN, la 1^{ère} personne du singulier s'écrit sans la *hamza* : باكتب. CODA recommande de ne jamais transcrire منكتب à la 1^{ère} personne du pluriel, même si c'est prononcé ainsi. Pour AraPI à l'inverse, on note منكتب.

Les prépositions

Les particules monolitères sont attachées : والكتاب - بالبيت

Les autres particules ne doivent pas être cliticisées : ما de négation, ainsi que ل + suffixe.

Par exemple, on écrira ما قلت لهاش (pas ماقلتلهاش).

Les démonstratifs

Les formes brèves comme *ha-* sont attachées dans CODA et AraPI, séparées par une espace dans BBN :

هالبيت - بهالبيت - عهالبيت vs ها البيت

Notation des chiffres et nombres en toutes lettres, jours de la semaine

Pour CODA et AraPI, les chiffres et les nombres écrits en toutes lettres sont conformes à l'écriture standard, quelle que soit la réalisation. Il en va de même pour les jours de la semaine. Pour BBN les nombres sont écrits en transcription dialectale (...طنعشر).

Les mots courants (*chose, eau, viens !*)

Pour CODA, on suit les règles du MSA dans les cas suivants :

Convention	Variantes non retenues
شيء	شي - إشي
ماء	ماي - مي - ماية
تعالوا - تعالوا	تع - تعي - تعوا

AraPI propose une liste ouverte de mots dialectaux et expressions très courantes.

5. Bibliographie complémentaire

La convention CODA appliquée à différents dialectes :

- HABASH, DIAB et RAMBOW, *Conventional Orthography for Dialectal Arabic (CODA): Principles and Guidelines - Egyptian Arabic* [En ligne]

<https://academiccommons.columbia.edu/doi/10.7916/D83X8562>

- HABASH, JARRAR, ALRIMAWI, AKRA, ZALMOUT, BARTOLOTTI et ARAR, *Palestinian Arabic Conventional Orthography Guidelines - Technical Report* [En ligne]

<http://www.jarrar.info/publications/HR15.pdf>

- SAADANE et HABASH, *A Conventional Orthography for Algerian Arabic* [En ligne]

<http://www.aclweb.org/anthology/W15-3208>

- TURKI, ADEL, DAOUDA et REGRAGUI, *A Conventional Orthography for Maghrebi Arabic* [En ligne]

https://www.researchgate.net/publication/311589181_A_Conventional_Orthography_for_Maghrebi_Arabic

- ZRIBI, BOUJELBANE, MASMOUDI, ELLOUZE, BELGUITH et HABASH, *A Conventional Orthography for Tunisian Arabic* [En ligne]

Annexe 5 : Tableau synoptique de la convention AraPI

Toutes les notations conventionnelles figurent en annexe 1 (phénomènes interactionnels) et annexe 2 (liste des notations conventionnelles).

Éléments conventionnels	Tier 1	Tier 2a	Tier 3	Tier 4
Transcription (type)	interactionnelle phonético-phonologique	morpho-syntaxique (référant arabe standard et standard dialectal)		traduction dans la langue du chercheur
Alphabet	API modifié		arabe (modifié) / latin	latin
Segmentation	en mots	conventionnelle morphémique par tiret (-)	en mots selon les règles de la graphie arabe	-
Interaction	ICOR modifiée	aucune notation sauf pauses, enchaînements et chevauchements		
Pronoms, déterminants, particules, prépositions	transcription phonétique	notation conventionnelle (référant arabe standard et standard dialectal)		-
Chiffres	transcription phonétique	notation conventionnelle (référant arabe standard et standard dialectal)		-
Morphologie nominale et verbale	transcription phonétique	notation conventionnelle (référant arabe standard et standard dialectal)		-
Expressions figées	transcription phonétique	codage 4 majuscules	notation conventionnelle	codage 4 majuscules
Mots ou syntagmes en langue étrangère	entre dièses, dans l'orthographe de la langue		entre dièses, transcription dans l'orthographe de la langue que les segments soient courts ou longs	entre dièses, traduits, mention de la langue par le code ISO entre doubles parenthèses
Mots étrangers entrés dans les usages	transcription phonétique	entre dièses, dans l'orthographe de la langue	entre dièses, transcription en caractères arabes	
Emprunts arabisés	entre dièses, transcription phonétique			
Régulateurs	transcription conventionnelle			
Noms propres	entre guillemets hauts orthographe officielle ou transcription phonétique avec majuscules	entre guillemets hauts orthographe officielle ou transcription standard avec majuscules	entre guillemets hauts orthographe officielle ou transcription en caractères arabes	orthographe officielle avec majuscules
Éléments dialectaux courants	transcription phonétique	notation conventionnelle en API modifié	notation conventionnelle en caractères arabes	-