

Introduction to Arabic Natural Language Processing (part 1)

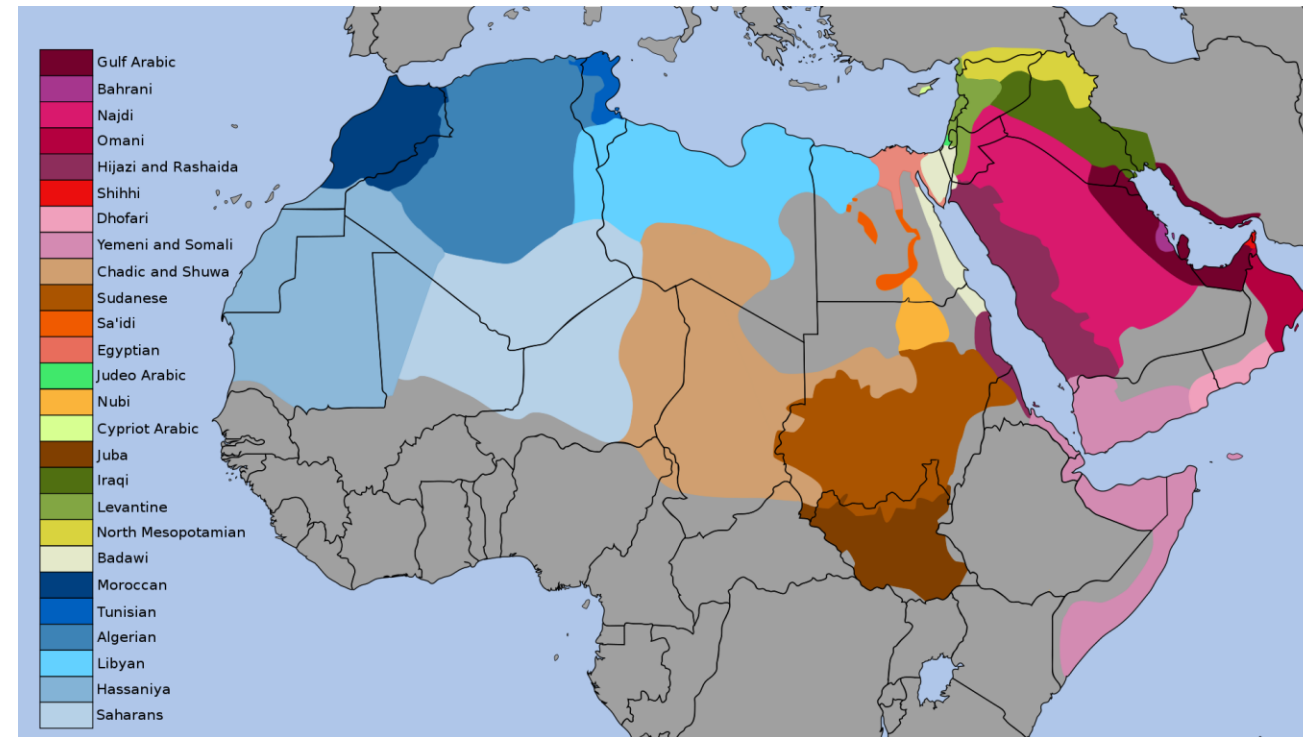
Kareem Darwish
 aixplain Inc.

المحاضرة الأولى: مقدمة



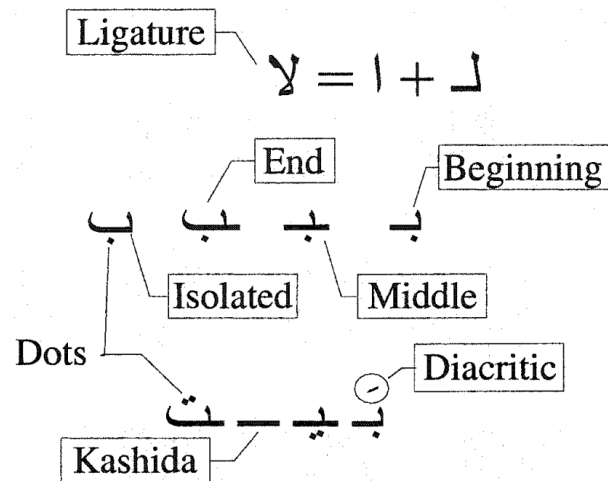
What is Arabic?

- Arabic is a Semitic language that is being spoken by more than 450 million people.
- Arabic was the language of science and culture for hundreds of years.
- Arabic script is being used by many languages (ex. Urdu, Farsi, Kurdi, Hausa, etc.), and used to the script for many more languages (ex. Turkish, Russian, etc.)
- Arabic has fundamental influence on many world languages (ex. Urdu, Maltese, Hebrew, Turkish, Farsi, Somali, Mahri, etc.)
- Arabic is a family of languages (ex. Arabic, Maltese, Mahri, Suqutri, Shehri, etc.) and dialects (ex. Maghrebi, Hassaniya, etc.).



How is Arabic Different – Orthography 1

- Orthography refers to how a language is written
- Arabic:
 - Is written from right-to-left
 - Uses connected characters that change shape depending on position
 - Is typically written without diacritics – diacritics are in fact part of words (ex. رَجُلٌ، رَجُل، رَجُل)



قَالَ بَلَىٰ وَلَٰكِن لِّيَطْمَئِنَّ قُلُوبِي قَالَ فَخُذْ أَرْبَعَةً مِّنَ الطَّيْرِ فَصُرْهُنَّ إِلَيْكَ ثُمَّ أَجْعَلْ عَلَىٰ كُلِّ جَبَلٍ مِّنْهُنَّ جُزْءًا ثُمَّ ادْعُهُنَّ يَأْتِينَكَ سَعْيًا وَاعْلَمْ أَنَّ اللَّهَ عَزِيزٌ حَكِيمٌ

How is Arabic Different – Orthography 2

- Arabic:
 - Has multiple scripts:
 - Othmani script
 - Maghrebi script
 - Arabizi
 - Has multiple encoding:
 - Buckwalter
 - UTF-8
 - Windows CP1256

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ
وَالسَّمَاءِ وَالطَّارِقِ ❶ وَمَا أَدْرَاكَ مَا الطَّارِقُ ❷ النَّجْمُ الثَّاقِبُ ❸ إِنَّ كُلَّ
نَفْسٍ لَّمَّا عَلَيْهَا حَافِظٌ ❹ فَلْيَنْظُرِ الْإِنْسَانُ مِمَّ خُلِقَ ❺ خُلِقَ مِنْ مَّاءٍ
دَافِقٍ ❻ يَخْرُجُ مِنْ بَيْنِ الصُّلْبِ وَالتَّرَائِبِ ❼ إِنَّهُ عَلَى رَجْعِهِ لَقَادِرٌ ❽
يَوْمَ تُبْلَى السَّرَائِرُ ❾ فَمَا لَهُ مِنْ قُوَّةٍ وَلَا نَاصِرٍ ❿ وَالسَّمَاءِ ذَاتِ الرَّجْعِ ❶❶
وَالْأَرْضِ ذَاتِ الصَّدْعِ ❶❷ إِنَّهُ لَقَوْلُ فَصْلٍ ❶❸ وَمَا هُوَ إِلَّا هَزْلٌ ❶❹ إِنَّهُمْ
يَكِيدُونَ كَيْدًا ❶❺ وَأَكِيدُ كَيْدًا ❶❻ فَهَلْ الْكَافِرِينَ أَهْمُ لَهُمْ رُؤُودًا ❶❼

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ
وَالسَّمَاءِ وَالطَّارِقِ ❶ وَمَا أَدْرَاكَ مَا الطَّارِقُ ❷ النَّجْمُ الثَّاقِبُ ❸ إِنَّ كُلَّ
نَفْسٍ لَّمَّا عَلَيْهَا حَافِظٌ ❹ فَلْيَنْظُرِ الْإِنْسَانُ مِمَّ خُلِقَ ❺ خُلِقَ مِنْ مَّاءٍ
دَافِقٍ ❻ يَخْرُجُ مِنْ بَيْنِ الصُّلْبِ وَالتَّرَائِبِ ❼ إِنَّهُ عَلَى رَجْعِهِ لَقَادِرٌ ❽
يَوْمَ تُبْلَى السَّرَائِرُ ❾ فَمَا لَهُ مِنْ قُوَّةٍ وَلَا نَاصِرٍ ❿ وَالسَّمَاءِ ذَاتِ الرَّجْعِ ❶❶
وَالْأَرْضِ ذَاتِ الصَّدْعِ ❶❷ إِنَّهُ لَقَوْلُ فَصْلٍ ❶❸ وَمَا هُوَ إِلَّا هَزْلٌ ❶❹ إِنَّهُمْ
يَكِيدُونَ كَيْدًا ❶❺ وَأَكِيدُ كَيْدًا ❶❻ فَهَلْ الْكَافِرِينَ أَهْمُ لَهُمْ رُؤُودًا ❶❼

How is Arabic Different – Orthography 2

- Arabic:
 - Has different writing conventions:
 - Ex. Terminal: ي (ex. علي) vs ى (ex. على)
 - Ex. Default diacritic: قَالَ vs قال

How is Arabic Different – Orthography 3

Arabic letters	ا	ب	ت	ث	ج	ح	خ	د	ذ	ر	ز	س	ش	ص	ض	ط	ظ	ع	غ	ف	ق	ك	ل	م	ن	هـ	و	ي	ى ^[4]
DIN 31635	' / ā			t̤	ǧ	ħ	ħ		d̤				š	ṣ	ḍ	ṭ	ẓ	ʿ	ġ								w / ū	y	ī
Buckwalter	A	b	t	v			x	d	*	r	z	s	\$					E	g	f	q	k	l	m	n	h	w	y	Y
Qalam	' / aa			th	j	H	kh		dh				sh	S	D	T	Z	`	gh									y	
BATR	A / aa			c			K		z'				x					E	g								w / uu	y	ii
IPA (MSA)	ʔ, aː	b	t	θ	dʒ	ħ	x	d	ð	r	z	s	ʃ	sˤ	dˤ	tˤ	ðˤ	ʕ	ɣ	f	q	k	l	m	n	h	w, uː	j, iː	

- Arabic text

- يُولَدُ جَمِيعُ النَّاسِ أَحْرَارًا مُتَسَاوِينَ فِي الْكَرَامَةِ وَالْحُقُوقِ.

- Buckwalter transliteration

- yuwladu jamiyEu {ln~aAsi >aHoraArFA mutasaAwiyna fiy {lokaraAmapi wa{loHuquwqi.

hamza

- lone hamza: '
 - hamza on alif: >
 - hamza below alif: <
 - hamza on wa: &
 - hamza on ya: }

alif

- madda on alif: |
- alif al-wasla: {
- dagger alif: `
- alif maqsura: Y

harakat

- fatha: a
- damma: u
- kasra: i
- fathatayn: F
- dammatayn: N
- kasratayn: K
- shadda: ~
- sukun: o

ta marbouta: p

How is Arabic Different – Morphology 1

- Morphology refers to how a word is constructed
- Arabic has a root-based morphological system:
 - Start with root (ex. كتب – ktb) – about 10,000 roots
 - Fit into *stem template*, ex.:
 - CCAC: كتاب – ktAb
 - CACC: كاتب – kAtb
 - mCCwb: مكتوب – mktwb
 - Add diacritics: كِتَاب – kitaAb, كُتِبَ – kut~aAb
 - Add prefixes and suffixes:
 - Ex. Prefixes: coordinating conjunctions (و، ف); determiners (ال); future (س); etc.
 - Ex. Suffixes: pronouns (هم، كم، هـ، ها، ك، كما، كن،) (ين، ون، ان، ي، و، ا) plural/dual markers (هن); etc.
 - Ex. وَكِتَابُهُنَّ – wakitaAbhun~a

* ك ت ب – (كَتَبَ) من بابِ نصر
و (كَتَابًا) أيضًا و (كِتَابَةً) . و (الكِتَابُ)
أيضًا الفَرْضُ والحُكْمُ والقَدَرُ. و (الكَاتِبُ)
عندَ العربِ العَالَمُ ومنه قوله تعالى :
« أَمْ عِنْدَهُمُ الْغَيْبُ فَهُمْ يَكْتُبُونَ »
و (الْكُتَّابُ) بالضمِّ والتشديدِ (الْكُتْبَةُ) .
و (الْكُتَّابُ) أيضًا و (المَكْتُبُ) واحدٌ^(٢)
والْجَمْعُ (الْكُتَاتِبُ) و (المَكَاتِبُ) .
و (الْكُتَيْبَةُ) الجَيْشُ . و (أَكْتُبُ) أي

How is Arabic Different – Morphology 2

- Arabic morphology is ambiguous:

- Ex. وليد:

- وليد – *waliyd*
- وليد – *wa+li+yad*

- Ex. ايمان:

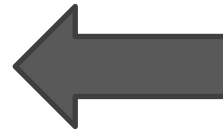
- ايمان -- *<iymaAn* – أم ن
- ايمان -- *>ayomaAn* – ي م ن
- ايمان -- *>ay~ima+An* – أي ن
- ايمان -- *>a+yamo>an* – م أن

وليد
ايمان

How is Arabic Different – Morphology 3

- One of the fundamental problems in Arabic is word segmentation:

ويكيبيديا هي موسوعة+ة رقمي+ة ، متعدد+ة
ال+لغ+ات ، حر+ة ال+محتوى . يستطيع أي
شخص ال+تحرير في+ها ب+دون تسجيل ،
و+يستطيع أي شخص ال+استفادة من
ال+محتوى ، و+استغلال+ه ب+هدف تجاري أو
غير+ه وفق+ال+ترخيص ال+موسوعة+ة .



ويكيبيديا هي موسوعة رقمية ، متعددة اللغات ، حرة
المحتوى . يستطيع أي شخص التحرير فيها بدون
تسجيل ، ويستطيع أي شخص الاستفادة من
المحتوى ، واستغلاله بهدف تجاري أو غيره وفقًا
لترخيص الموسوعة .

How is Arabic Different – Morphology 3

- Segmentation is important for many downstream applications such as:
 - Information retrieval (search)
 - كتاب → والكتاب، الكتاب، كتابهم، وكتابهن
 - Lemmatization (الرجوع لأصل الكلمة):
 - يد;وليد → وليد
 - رجل → الرجال
 - Syntactic processing:
 - و+ل+يد → وليد: noun; و+ل+يد: coordinating_conjunction + preposition + noun
 - Diacritic recovery:
 - وَلِيد، وَ ل+يَد → وليد

ويكيبيديا هي موسوعة رقمية ، متعددة اللغات ، حرة المحتوى . يستطيع أي شخص التحرير في ها
بدون تسجيل ، ويستطيع أي شخص ال استفادة من المحتوى ، واستغلال ها بهدف تجاري أو
غير ها وفق ال ترخيص الموسوعة .

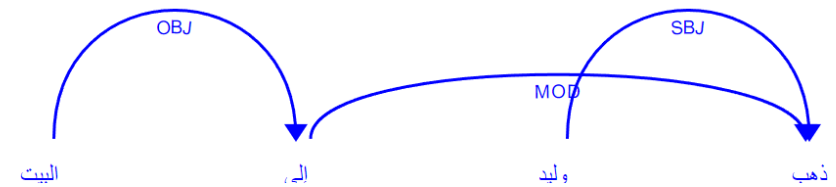
ويكيبيديا هي موسوعة رقمية ، متعددة لغة ، حرة المحتوى . استطاع أي شخص تحرير في دون تسجيل ، استطاع أي
شخص استفادة من محتوى ، استغلال هدف تجاري أو غير وفق ترخيص موسوعة .

NOUN PRON NOUN+NSUFF ADJ+NSUFF PUNC ADJ+NSUFF
DET+NOUN+NSUFF PUNC ADJ+NSUFF DET+NOUN PUNC V NOUN NOUN
DET+NOUN PREP PRON PREP NOUN NOUN PUNC CONJ V NOUN NOUN
DET+NOUN+NSUFF PREP DET+NOUN PUNC CONJ NOUN PRON PREP
NOUN ADJ CONJ NOUN PRON NOUN CASE PREP NOUN
DET+NOUN+NSUFF PUNC

ويكيبيديا هي موسوعة رقمية ، متعددة اللغات ، حرة المحتوى . يستطيع أي شخص التحرير فيها بدون تسجيل ،
ويستطيع أي شخص الاستفادة من المحتوى ، واستغلاله بهدف تجاري أو غيره وفقاً لترخيص الموسوعة .

How is Arabic Different – Morphology 4

- Segmentation is important for many downstream applications such as:
 - Parsing (إعراب):
 - ذهب وليد إلى البيت



Dependency Relation	HEAD ID	FORM	ID
---	0	ذهب	1
SBJ	1	وليد	2
MOD	1	إلى	3
OBJ	3	البيت	4

How is Arabic Different – Syntax 1

- Syntax refers to how a sentence is constructed
- Arabic is generally free form:
 - Verb-Subject-Object (VSO): أكل الأسد الأرنب
 - Subject-Verb-Object (SVO): الأسد أكل الأرنب
 - Verb-Object-Subject (VOS): أكل الأرنب الأسد
 - Subject-Predicate (SP - جملة اسمية): الأسد مفترس
- Two central problems are associated with syntax, namely:

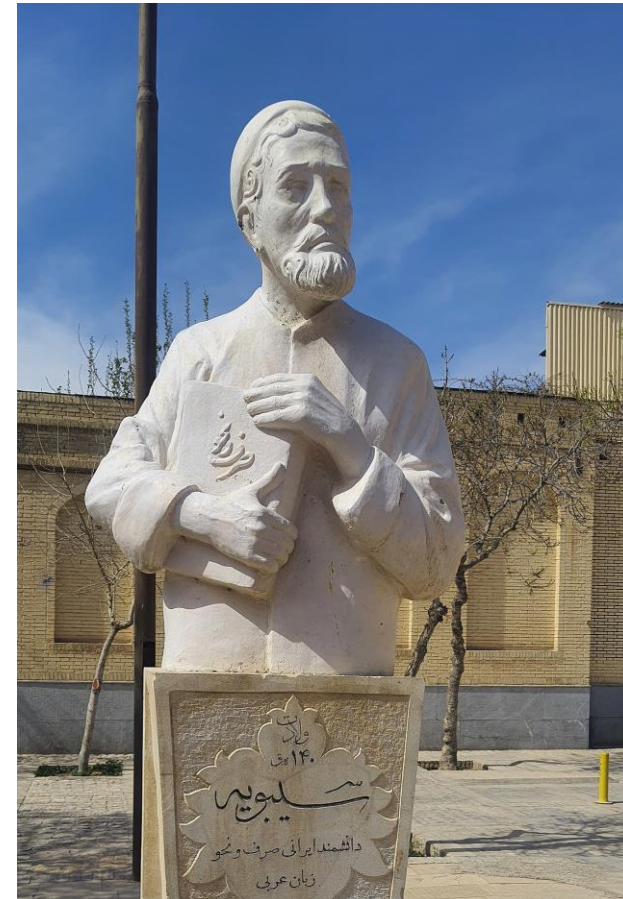
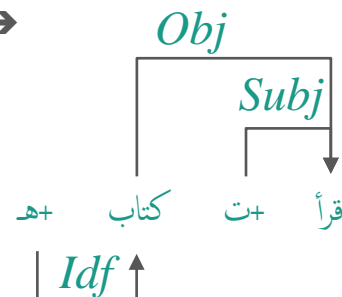
- Part of speech (POS) tagging:

■ Ex. قرأت كتابه →

Pronoun	هـ+	Noun	كتاب	Pronoun	ت	Noun	قرأ
---------	-----	------	------	---------	---	------	-----

- Parsing:

■ Ex. قرأت كتابه →



Diacritization Convention

Diacrtization is NOT Always Easy

فكاهة: من نظر إلى أخيه نظرة ...

مثلث قطرب!

خلا أخ لك فلا تخله ... من لك يوما بأخيك كله!

إن الدنيا أقبلت بصرم وأدبرت حذاء، ولم يبق منها إلا كصابة الإناء يشربها أحدكم



Diacrtization is NOT Always Easy

بدا فحيا بالسَّلام ~ رمى عذولي بالسِّلام

أشار نحوي بالسُّلام ~ بكفه المخضَّب

بالفتح لفظُ المبتدي ~ والكسرِ صخرُ الجَلَمَدي

والضمِ عِرْقُ في اليدي ~ قد جاء في قول النبي

Diacrtization is NOT Always Deterministic

سَاعَدَ الرَّجُلُ الْوَلَدَ vs سَاعَدَ الرَّجَلَ الْوَلَدُ
أَطْعَمْتُ الْهَرَّةَ vs أُطْعِمْتُ الْهَرَّةَ vs أَطْعَمْتُ الْهَرَّةَ



Which form is correct?

بَرَزَ الثَّعْلَبُ يَوْمًا فِي ثِيَابِ الْعَارِفِينَ ... وَمَشَى فِي الْأَرْضِ يَهْدِي وَيَسُبُّ الْمَاكِرِينَ
بَرَزَ الثَّعْلَبُ يَوْمًا فِي ثِيَابِ الْعَارِفِينَ ... وَمَشَى فِي الْأَرْضِ يَهْدِي وَيَسُبُّ الْمَاكِرِينَ

قَالَتْ: دَعِينِي وَهْزَالِي وَالزَّمَنُ ... وَأَجِيبِي صَاحِبَ السِّكِينِ يَا ذَاتَ الثَّمَنِ
قَالَتْ: دَعِينِي وَهْزَالِي وَالزَّمَنُ ... وَأَجِيبِي صَاحِبَ السِّكِينِ يَا ذَاتَ الثَّمَنِ

Which form is correct?

ثَوَى طَاهِرُ الْأُرْدَانِ . لَمْ تَبَقْ بُقْعَةٌ ... غَدَاةَ ثَوَى إِلَّا اشْتَهَتْ أَنَّهَا قَبْرُ
ثَوَى طَاهِرِ الْأُرْدَانِ . لَمْ تَبَقْ بُقْعَةٌ ... غَدَاةَ ثَوَى إِلَّا اشْتَهَتْ أَنَّهَا قَبْرُ

اشْتَقْتُ لِرُؤْيَيْهِ

اِشْتَقْتُ لِرُؤْيَيْهِ

Conventions R

Important 4 Consistency

Should we explicitly put default diacritics? Ex. قَالَ vs قَالْ

Should we put explicit *sukun* on letters? Ex. مِثْلُ vs. مِثْلٌ

Should we put diacritics on different forms of *hamza*? أَنْتَ vs. أَنْتَ

Should we put diacritics on *hamzat-ul-wash*? اِضْرِبْ vs. اضْرِبْ

Should we put *sukun* on *lam* in *Al Alqamariyya*? الْقَرْيَةُ vs الْقَرْيَةُ

Formulating a Convention for Gulf

Let's look at some examples:

لوی راسه بسرعة وقام يصيح: ش تقول إنت ش تقول!

ابتسمت بخبث ورمست: شو سالفتيج

اتصل وخل الباقي على الله

وربع هزاع صوبه ومط الفون منو وسكره بوجهه



Dialectal Diacritization NEEDS Conventions

In Gulf dialect, which diacritized form is correct:

تَقُولُ vs تَقُولُ vs تَقُولُ

Can we start a word with *sukun*?

Can we have two letters with *sukun* back-to-back? Ex. عِنْدَهَا

Are we going to put ANY case endings (حركة إعراب)? Ex. عَبْدَ اللَّهِ vs عَبْدُ اللَّهِ

How do we diacritize the coordinating conjunction و?

Sample Convention

- الحرف الأول يجب تسكينه إذا كان يسكن مع إضافة حرف الواو قبله. مثال: كلمة "يَقُول" يسكن أولها لأنه ما إضافة "و" قبلها تشكل كالاتي: "وَيَقُول". إذا حرك مع وضع "و" قبل، إذا توضع الحركة المناسبة. مثال: "بِتَقُول"
- عند الشك بين الكسر أو تشكيل آخر، يغلب الكسر. مثال: الواو في "وَالْحِينَ"
- تكسر "و" الإضافة في أغلب الأحيان (مثال: "وَهِيَ") إلا إذا كان هناك حركة أخرى ظاهرة
- التقاء الساكنين مسموح به في حدود ضيقة. مثل: "بَعْد" f
- يترك الحرف الأخير دون تشكيل حيث إن اللهجات ليس بها حركة إعراب، ولكنه في بعض الأحيان يحتاج إلى شدة (مثل: "رَدّ")

Sample Convention

- عند التشكيل بالفتح يوضع التنوين قبل الألف. مثال: "طَبْعًا"
 - عند المد بالألف أو الواو أو الياء، يجب وضع حركة على الحرف التي قبلها. مثال: "قَالَ"
 - يجب تشكيل جميع الحروف إلا حروف المد وهمزة الوصل: مثال: "اسْتَحْت"
 - يجب التشكيل حسب طريقة النطق في اللهجة
 - تسكن واو الجماعة في أغلب الأحيان (مثال: "شَافُوا"), إلا إذا استلزم وضع ضمة على الحرف التي قبلها (مثال: رَدُّوا)
-

Sample Convention

هذه بعض الجمل المشكّلة:

• هَزَّاعٌ : الرِّيمُ قَالَتْ كُلُّ شَيْءٍ

• مَانِعٌ : قُولِي لَا خَلِّيتِ مَنِّجَ إِنْ شَاءَ اللَّهُ

• وَقَفَّني صُوتَ عَمِّي خَالِدٌ

• نَاصِرٌ : إِنِّي مَا لَجُ خَصَّ جَزُوي سِيرِي وَفِكِّي

• الرِّيمُ : عَبَّالِي ^ ^

• مِيرَةٌ : يَعْنِي عُقْبُ سَاعَةٍ

• فَزَّ خَلِيفَةً : هَلَا جُدُّوهُ ..