

Authentic Audio Resources for “Kitchen Table Arabic” and Code-Switching Patterns

Kitchen Table Arabic refers to the informal, everyday dialect spoken in homes and streets, as opposed to formal textbook Arabic ¹. For Syrian and Lebanese refugees, this colloquial speech often carries heavy emotional content – including **somatic idioms of distress** like “*my chest is heavy*” or “*my liver is burning*” to express grief ² – and frequently mixes in French or English words (code-switching) ³. Below, we identify sources of **authentic audio** in this dialect (especially from refugees in trauma contexts), and outline **common code-switching patterns** in Lebanese Arabic, along with corpora for training speech/text models on dialectal Arabic.

Real-Life Syrian Refugee Audio Samples (Emotion & Code-Switching)

- **Community Oral Histories (Public Archives):** The *Syrian Voices* project at Oxford recorded over 50 in-depth interviews with Syrian refugees, capturing personal narratives of war, displacement, and resettlement ⁴ ⁵. These filmed and audio-recorded testimonies provide genuine “*kitchen table*” dialogue – often emotionally charged accounts in Syrian Arabic. (Contact: info@syrianvoices.org.uk for access to the archive ⁶.) Another example is the “**A Million Stories**” digital library (an EU-funded archive of refugee experiences), which includes audio recordings of Syrian refugees sharing their stories in their native dialect. These oral history resources feature authentic voices from the field, though access may require reaching out to project organizers or archives.
- **Humanitarian and Journalistic Field Recordings:** Photojournalist Tanya Habjouqa’s multimedia story for *TIME* (“The Messages That Hold Refugee Families Together”) uniquely used WhatsApp **voice messages** from Syrian refugees to convey emotional reality ⁷. In one case, a mother in a refugee camp plays an *audio message* of her husband (separated in Germany) singing a lullaby to their child – an intimate example of real refugee audio capturing *love, hope and sorrow* ⁷. Habjouqa collected dozens of such voice notes and wove them into a short documentary video ⁷. This kind of journalistic project provides **authentic dialectal speech** (often **Levantine Arabic with some English** phrases) in highly emotional contexts. Similarly, various NGO and UN agencies have recorded refugee testimonies for advocacy: e.g. UNHCR and UNICEF have published video interviews (with original audio) of Syrian women and children describing trauma, which researchers can sometimes access via the UN Audiovisual Library or organization websites. These sources are often publicly viewable (e.g. on YouTube or UN media sites) and feature spontaneous speech, occasionally code-switching when refugees use an English word or local French term.
- **Academic Interview Datasets:** Researchers studying refugee mental health have also **audio-recorded field interviews** in Arabic. For example, a 2020 psychosocial study in Jordan conducted 24 *in-depth interviews* with Syrian refugee women (in Arabic, with translation/transcription) about war trauma and displacement ⁸. The recordings (used for qualitative analysis in *Journal of*

Psychosomatic Research) captured emotional narratives including descriptions of physical symptoms of trauma ⁹. While such datasets are typically not open public files (due to privacy/ethics), the **transcripts** often preserve **somatic idioms** and emotional expressions – e.g. women describing how “*war and displacement made their bodies and hearts heavy*” (somatization of mental anguish) ⁹. Interested parties could contact the authors (N. Rizkalla et al.) for potential access or locate supplementary material if provided. Academic projects focusing on **trauma healing** in MENA cultures sometimes compile a “**Dictionary of Pain**” mapping these idioms ², which can inform language models even if the raw audio is restricted.

- **Licensed Speech Corpora (Levantine Dialects):** For training speech-to-text on dialectal Arabic, several large corpora include *spontaneous Levantine speech*. Notably, the **LDC (Linguistic Data Consortium)** catalogs: **Fisher Levantine Arabic Conversational Speech** – 45 hours of unscripted phone calls in Levantine dialect (Syrian, Lebanese, Jordanian, etc.) ¹⁰ – and the earlier **CallHome Syrian Arabic** telephone corpus (colloquial family calls). These contain everyday “kitchen table” conversations (topics like family, daily life) with authentic dialect and some natural emotion. While not specific to refugee trauma, such data provides a **baseline for colloquial speech**. Access requires LDC membership or licensing ¹¹ ¹⁰, and transcripts are available ¹².
- Of special interest, **BBN/AUB Babylon Levantine Corpus** (LDC2005S08) was collected to build a **refugee/medical speech translator**. It has ~60 hours of Levantine colloquial speech, recorded via prompts that simulate refugee interactions in medical or aid settings ¹³. For example, subjects (in Lebanon and the U.S.) responded to scenarios like medical triage or aid interviews, producing utterances useful for a “handheld” translator system ¹³. This means the content likely includes phrases about health, needs, locations – potentially mirroring refugee speech (“I need water,” “My child is sick,” “I lost my home,” etc.) in dialect. The audio and matching transcripts (76k utterances) are available from LDC ¹⁴ ¹⁵ under license. This corpus can be highly relevant for training ASR or dialogue systems to understand *informal Syrian/Lebanese Arabic in humanitarian contexts*.
- **Other Sources:** There are many **personal storytelling podcasts and radio pieces** featuring Syrian refugees. For instance, SBS Australia’s Arabic program recorded Syrian women in a Lebanese camp telling their stories (in Syrian Arabic) for a special broadcast – capturing their voices and even code-switched bits when they use a borrowed word. Community initiatives like **Refugee Voices** (by Translators Without Borders) have collected audio interviews to improve translation services ¹⁶ – volunteers transcribe refugees’ spoken narratives, which indicates such audio exists (often stored by NGOs). While these may not be packaged as formal datasets, they are valuable **real speech samples**. Reaching out to humanitarian organizations or media outlets can often yield contacts or links to archives of these audio files (sometimes available on SoundCloud, YouTube, or institutional sites).

Tip: When using these audio resources, check licensing. Projects like Syrian Voices are for research/educational use (contact required), and LDC corpora require purchase or membership (download via LDC catalog links ¹²). Journalistic content (YouTube videos, etc.) is usually public – one can directly download audio from those with permission or via the publisher’s site. Each resource above prioritizes **authenticity** – real voices of refugees speaking informal Arabic – which is crucial for training empathetic language models.

Code-Switching in Lebanese Arabic – Patterns and Corpora

Lebanese Arabic is famously **mixed with French and English** in daily conversation, a legacy of Lebanon's multilingual culture ¹⁷. Lebanese speakers (and Syrians in Lebanese diaspora) often alternate between languages even within a single sentence (so-called *intra-sentential* switching) ¹⁷. This code-switching is so natural that even short phrases can contain three languages. For example, a typical greeting might be: "Hi! Kifak, ça va?" – English "hi," Arabic "how are you," and French "how's it going?" all in one breath ¹⁸. Below are **common code-switching patterns** and relevant linguistic resources:

1. Tri-Lingual Greetings and Farewells: It's common to open or close conversations by switching languages for each part. The greeting above ¹⁸ is one popular example. Similarly, one might say "Habibi, see you, bisous!" mixing Arabic (habibi, "dear"), English ("see you [soon]"), and French ("kisses") ¹⁹. Many Lebanese use "Yalla, bye" (Arabic "let's go/okay" + English "bye") to say goodbye – an expression now iconic of Lebanese code-mix. These routine social phrases seamlessly blend languages.

2. Politeness and Exclamations in French/English: French is often used for courtesy terms or interjections in otherwise Arabic sentences. "**Merci!**" (thank you) and "**bonjour/bonsoir**" (hello/good evening) are nearly as common as their Arabic equivalents ²⁰. For instance, a Lebanese diner might call a waiter with "**Maître! El-fetoura, please.**" – here "Maître" (French honorific for a headwaiter) and "please" (English) frame the Arabic "the bill". Likewise, you'll hear "**Pardon,**" "**Excuse me,**" or "**Sorry**" in English inserted into Arabic dialogue for politeness. These borrowed niceties have become part of the Lebanese dialect repertoire ¹⁷ ¹⁸.

3. Common French Loanwords for Everyday Nouns: Lebanese Arabic has absorbed **hundreds of French nouns** over decades. Speakers often default to the French term for certain objects, especially if it has no easy Classical Arabic equivalent in casual speech. Some examples frequently cited by Lebanese themselves include: "*toilette*" (bathroom), "*chapeau*" (hat), "*jaquette*" (jacket), "*chauffeur*" (driver), "*abajour*" (lampshade), "*piscine*" (swimming pool), etc. ²¹ ²². Many of these are so integrated that they're pluralized or pronounced with Arabic style. English words have entered in a similar way for modern items – e.g. "**phone**", "**TV**", "**computer**" – though French dominated older generations. Any language model for Lebanese speech must handle these embedded foreign terms as part of the lexicon.

4. Switching for Technical and Academic Terms: Due to Lebanon's educational system (where science and math are taught in French or English), educated Lebanese switch to the foreign term when discussing technical topics ²³. A study found that discussions on medicine or engineering naturally produced more code-switching – the speaker would use the exact French/English term rather than finding an Arabic substitute ²³. For example, someone might say "عندی diabetes, بحاجة insulin" ("I have diabetes and need insulin"), inserting the English medical terms. **Keeping the foreign term intact** is often deliberate: it ensures clarity (everyone knows the term as taught) and signals the speaker's education. This pattern was identified as a primary function of code-switching in Lebanon ²⁴. It also extends to university jargon, IT and finance (e.g. "لأجل upgrade لـsoftware" – "I did an upgrade to the software").

5. Emphasis and "Showing Off" Knowledge: Sometimes switching language can serve to emphasize a point or to signal cosmopolitan identity. Lebanese speakers might quote a saying in French for effect, or use an English idiom to be humorous. The 2002 AUC thesis on Lebanese code-switching noted that one motivation was "*showing off knowledge of [foreign] languages*" in social contexts ²⁵. For example, dropping a clever English one-liner in an Arabic conversation, or using a French proverb, can index sophistication or

worldliness. This pattern is more about **socio-pragmatic function** than necessity – it's a stylistic choice common especially among younger urban Lebanese who are proudly trilingual.

6. Taboo or Emotional Topics in Another Language: Interestingly, speakers may switch to a foreign language when discussing very personal or culturally taboo subjects. Research in code-switching notes that talking about sex or deep feelings in a second language can feel “one step removed,” making it psychologically easier ²⁶. In Lebanon, for instance, a person might use English for a sexual term or French for a romantic endearment, whereas in pure Arabic it might feel too direct or embarrassing ²⁶. Similarly, someone grieving might slip into English for words like “depressed” or “anxiety” if those carry less stigma in translation. This cross-language shift helps speakers navigate cultural sensitivities while still communicating the idea.

7. Filler Words and Tags: Code-switching also occurs at the level of filler words (e.g. “you know...”, “like...”) and conversational tags. Lebanese Arabic speakers often say “OK,” “yeah,” or “I mean” within Arabic sentences. Even the Arabic filler “ya’ni” (“meaning/like”) can be followed by an English clause, and vice versa. Small words like “basically,” “so,” “anyway” are sprinkled into the flow. Tag questions are another example: “... mish heik? Right?” where “right?” in English confirms a statement made in Arabic. These micro-level switches are pervasive and often unconsciously used.

8. Numerals and Dates: In casual settings, Lebanese switch languages for numbers depending on context. French numerals (especially 70, 80, 90, which are long words in Arabic) are common – e.g. saying a phone number or price partly in French. A person might say “taljée b- 70 dollar” (francophone *soixante-dix* for 70), or give the year “deux mille vingt-trois” (2023) in French during an Arabic conversation. Younger speakers increasingly use English numbers too (in tech or business contexts). This numeric code-switching is practical and shaped by whichever system (French or English) the speaker finds quicker.

9. Arabic English Mix in Diaspora: Among Syrian refugees in Lebanon or those resettled in Western countries, you’ll hear hybrids like “شكراً kteer, thank you” or “يُوْهْ rooh home”. These are natural outcomes of picking up new languages. In fact, one analysis of a pan-Arab reality TV show (*Top Chef MENA*) found that **Lebanese participants code-switched far more than Syrians**, reflecting Lebanon’s multilingual norm ²³. (In that study, Lebanese contestants produced 58 switches vs. 14 by Syrian contestants in similar interactions, highlighting the contrast ²³.) Refugees who spent time in Lebanon might adopt Lebanese-style French insertions (like “merci kteer”) and those in English-speaking countries often mix English in daily Arabic (e.g. “I applied for **asylum**” with “asylum” in English). These mixed speech patterns should be considered in any corpus for model training.

Linguistic Resources & Corpora: To train or evaluate speech/text models on such code-switched dialectal Arabic, a combination of datasets is useful:

- **Annotated Code-Switch Corpora:** While Lebanese-specific corpora are rare, the phenomenon is covered in general Arabic CS datasets. The **QATAR Arabic Speech Corpus (QASR)** includes ~2,000 hours of broadcast audio (mostly Modern Standard Arabic, but ~30% dialectal) with some code-switching segments ²⁷ ²⁸. A subset called **QASR.CS** (5.9 hours) was extracted to test ASR on Arabic-English switching ²⁸. It contains mostly English insertions (and very few French) within Arabic news speech. Another is **ESCPWA-CS** (2.8 hours of recordings from a 2019 UN ESCWA meeting) which has extensive *intra-sentential* switching between Arabic and English ²⁹. In ESCWA-CS, about **35% of the speech segments include code-switching**, with ~22% of the content being English or French words

²⁹. (Notably, when North African delegates spoke, the switches were often Arabic-French ³⁰.) These corpora are available through the ArabicSpeech initiative (QCRI) – see arabicspeech.org – typically for research use. They provide **audio + transcripts with language tags**, ideal for testing speech recognizers on mixed language input.

- *North African Code-Switch Data:* Much research on Arabic code-switching has focused on Maghreb dialects (due to prevalent French mixing in Algeria, Morocco, etc.). For example, the **CAFE Corpus** (Codeswitching in Algerian French-English) provides ~37 hours of Algerian Arabic spontaneous speech with French switches ³¹. While not Lebanese, these datasets (available via ELRA or as described in papers) can help develop models for **Arabic-French code-mixing**, since the linguistic behavior is similar. They typically include audio and transcriptions with word-level language labels ³¹. *Insight:* A model trained on Algerian-French code-switch might generalize to Lebanese-French input (with some fine-tuning), given overlap in borrowed vocabulary.
- *Multi-Dialect Speech Collections:* Very recently (2024), a community-released corpus called **Casablanca** compiled **48 hours of fully transcribed speech** across 8 Arabic dialects (including Levantine), with explicit labeling of code-switched words ³² ³³. It tags English and French insertions in the transcripts (and even provides the foreign words in both Latin and Arabic script) ³³. Casablanca is positioned as the largest supervised Arabic dialect speech dataset, and it aims to fill exactly the gap of dialectal and code-mixed ASR training data ³⁴ ³². Once publicly released (likely via an open repository or by request to the authors ³⁵), this corpus will be highly suitable for training speech models that handle *Arabic dialect variance plus code-switching*. It can be combined with older corpora (e.g. MGB Challenge data for Egyptian, etc.) to improve robustness.
- *Text Corpora for Code-Switch Analysis:* In addition to audio, there are textual datasets of Arabizi (Arabic written in Latin script, often mixed with English/French). While not directly audio, they reveal common code-switch patterns in informal communication. For instance, the **Lebanese SMS/ Facebook Corpus** by Hamadeh (2017) studied youth messages with Arabic-English mix ³⁶. Such corpora list frequent foreign words and switching contexts in Lebanese digital communication. For training *language understanding*, these can augment speech data – e.g. providing examples of how a sentence might toggle between scripts or languages.

Download Links / Access: For *public datasets* like the QCRI corpora (ESCWA-CS, QASR-CS) and Casablanca, keep an eye on official websites or papers: the QCRI ones can be requested via their site (previous MGB Challenge data was distributed to researchers), and Casablanca's authors (Talafha et al. 2024) may release it on a platform like HuggingFace or by email request (contact info in the paper ³⁵). The LDC corpora (Fisher, CallHome, BBN/AUB) can be obtained from the LDC Catalog ¹⁴ ¹⁵ – we provided the catalog IDs for reference. If you don't have LDC access, some subsets (e.g. CallHome Levantine) have been used in ASR competitions and may have shorter samples available in publications. For the oral history and NGO recordings, usage is typically via the organization's site or on YouTube (e.g. *Syrian Voices* has a documentary film online, and *A Million Stories* shares excerpts on its website). We recommend reaching out to those organizations for raw audio files if needed for research.

By leveraging the above **authentic audio sources** and **linguistic studies**, one can gather a rich collection of *Kitchen Table Arabic* speech. These include the emotive, culturally nuanced expressions of Syrian refugees and the fluid code-switching idiom of Lebanese-influenced Arabic. Together, they will not only help train

robust speech-to-text models for dialectal Arabic, but also ensure the models truly “**listen with cultural empathy**” ³⁷ – recognizing phrases of pain and comfort in the speaker’s own tongue.

Sources:

- Maamouri, M. et al. *Fisher Levantine Arabic Conversational Telephone Speech*, LDC2007S02 – 279 calls (45 hours) of Jordanian, Lebanese, Palestinian, Syrian colloquial dialogues ¹⁰ (LDC catalog).
- BBN/AUB *DARPA Babylon Levantine Arabic Speech and Transcripts* – 60.6 hours of refugee/medical domain colloquial speech (Levantine dialect), audio + transcripts ¹³ (LDC2005S08).
- Rizkalla, N. et al. (2020). “*Women in refuge: Syrian women voicing health sequelae...*” – Qualitative interviews (audio-recorded in Arabic) with Syrian refugee women about war trauma ⁸ ⁹.
- **Syrian Voices** (2021–22, Oxford University) – 50+ filmed interviews of Syrian refugees in the UK, preserving authentic dialect and personal stories ⁴ ³⁸.
- Habjouqa, T. (2016). *Time Magazine* – “*The Messages That Hold Refugee Families Together*” – multimedia piece with WhatsApp audio notes of Syrian refugees (lullabies, voice messages) expressing hope and sorrow ⁷.
- Chowdhury, S.A. et al. (2021). *Interspeech* – Arabic code-switching ASR using **ESCWA-CS** (2.8h UN meeting) & **QASR-CS** (5.9h news) corpora ²⁹ ²⁸ (available via QCRI/ArabicSpeech).
- **Casablanca** Corpus (Talafha et al. 2024) – 48h multidialect Arabic speech dataset with **code-switch tags** for English/French content ³³ (contact authors for access).
- El-Samaty, M. (2002). *Thesis, AUC – Code-switching among educated Lebanese*: Identified functions and patterns (e.g. 62% of foreign words were code-switches; women use more French, men more English; switching often keeps technical terms in French/English) ²³.
- Transpremium Blog – “*Lebanon: The Ultimate Code-Switching Country*” (2016) – informal examples of Lebanese mixing (e.g. “*Hi! Kifak, ça va?*”, “*Habibi see you, bisous*”) ¹⁸ ¹⁹ illustrating everyday trilingual speech.
- Gardner-Chloros, P. (2009). *Code-Switching* (referenced in Caramel film study) – notes that Arabs may switch language for taboo or emotional topics to ease discussion ²⁶.
- Reddit r/Lebanon thread – “*French-Arabic words used in Lebanon*” – community-sourced list of common French loanwords in Lebanese Arabic (e.g. bonjour, merci, toilette, chapeau, jaquette, etc.) ²¹. (Anecdotal but insightful).

¹ ² ³ ³⁷ 1770047563-3789-glossary-of-ter.pdf
file:///file-X6Pno5QrmYey2HkukYepHu

⁴ ⁵ ⁶ ³⁸ Syrian Voices | TORCH | The Oxford Research Centre in the Humanities
<https://www.torch.ox.ac.uk/syrian-voices>

⁷ Syrian Refugees and WhatsApp: Hear Their Messages | TIME
<https://time.com/4272666/refugees-stories-whatsapp/>

⁸ ⁹ Women in refuge: Syrian women voicing health sequelae due to war traumatic experiences and displacement challenges - ScienceDirect
<https://www.sciencedirect.com/science/article/abs/pii/S002239991930594X>

¹⁰ ¹¹ ¹² Fisher Levantine Arabic Conversational Telephone Speech - Linguistic Data Consortium
<https://catalog.ldc.upenn.edu/LDC2007S02>

- 13 14 15** BBN/AUB DARPA Babylon Levantine Arabic Speech and Transcripts - Linguistic Data Consortium
<https://catalog.ldc.upenn.edu/LDC2005S08>
- 16** Looking for a Arabic translator for a Refugee project
<https://community.translatorswb.org/t/looking-for-a-arabic-translator-for-a-refugee-project/57619>
- 17 18 19 20** Lebanon: The Ultimate Code-Switching Country
<https://transpremium.com/lebanon-the-ultimate-code-switching-country/>
- 21 22** What are some French-Arabic words used in Lebanon? : r/lebanon
https://www.reddit.com/r/lebanon/comments/1j90hdv/what_are_some_frencharabic_words_used_in_lebanon/
- 23 24 25** "A study of code-switching among educated Lebanese as reflected in tele" by Mona El Sayed El Samaty
https://fount.aucegypt.edu/retro_etds/1613/
- 26** diva-portal.org
<https://www.diva-portal.org/smash/get/diva2:1390863/FULLTEXT01.pdf>
- 27 28 29 30** Towards One Model to Rule All: Multilingual Strategy for Dialectal Code-Switching Arabic ASR
https://www.isca-archive.org/interspeech_2021/chowdhury21_interspeech.pdf
- 31** Spontaneous code-switching speech dataset in Algerian dialect ...
<https://pmc.ncbi.nlm.nih.gov/articles/PMC12648582/>
- 32 33 34 35** Casablanca: Data and Models for Multidialectal Arabic Speech Recognition
<https://arxiv.org/html/2410.04527v1>
- 36** [PDF] Gender and linguistic background in SMS code-switching by ...
https://www.intercultural.urv.cat/media/upload/domain_317/arxius/TP5/08-Bassan.pdf