



HAL
open science

A Tale of Two Phonologies: Problems in Armenian ASR

Samuel Chakmakjian

► **To cite this version:**

Samuel Chakmakjian. A Tale of Two Phonologies: Problems in Armenian ASR. Journée d'apprentissage des représentations de la parole et du langage (GdR TAL), Apr 2022, Grenoble, France. . hal-04059138

HAL Id: hal-04059138

<https://hal.science/hal-04059138v1>

Submitted on 5 Apr 2023

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.

A Tale of Two Phonologies: Problems in Armenian ASR

INTRODUCTION

Armenian is a traditionally under-resourced language, which has seen a recent uptick in interest in the development of its tools and presence in the digital domain. Some of this recent interest has centered around the development of Automatic Speech Recognition (ASR) technologies. However, the language boasts two standard variants which diverge on several levels, the most salient of which being on the phonetic-phonological level. Despite this divergence several social factors render a unified ASR system preferable and practical. These factors include :

- a high level of mutual-intelligibility
- intense social overlap of the speech communities
- increased contact between the two; bi-variant households and speakers

THE PROBLEM

In terms of ASR, the most problematic feature of the divergence between Standard Western Armenian (SWA) and Standard Eastern Armenian (SEA), is that of the phonologies, and more specifically, the voicing systems.

SWA (Western Armenian)

- Traditionally: Anatolia, (post-)Ottoman zone
- Currently: Diasporan communities in Middle East, Europe and North America
- 30 phonemes : 24 consonants, 6 vowels
- **Two-way voicing system**
- Voiced; voiceless aspirated

	Bil.	L-d.	Alv.	P.-Al.	Pal.	Vel.	Uv.	Glott.
Nas.		m		n				
Plos.		p ^h b		t ^h d				k ^h g
Affr.				ts dz				tʃ dʒ
Fric.		f v		s z		ʃ ʒ		χ ʁ h
Appr.				l			j	
Tap				r				
Trill								

SEA (Eastern Armenian)

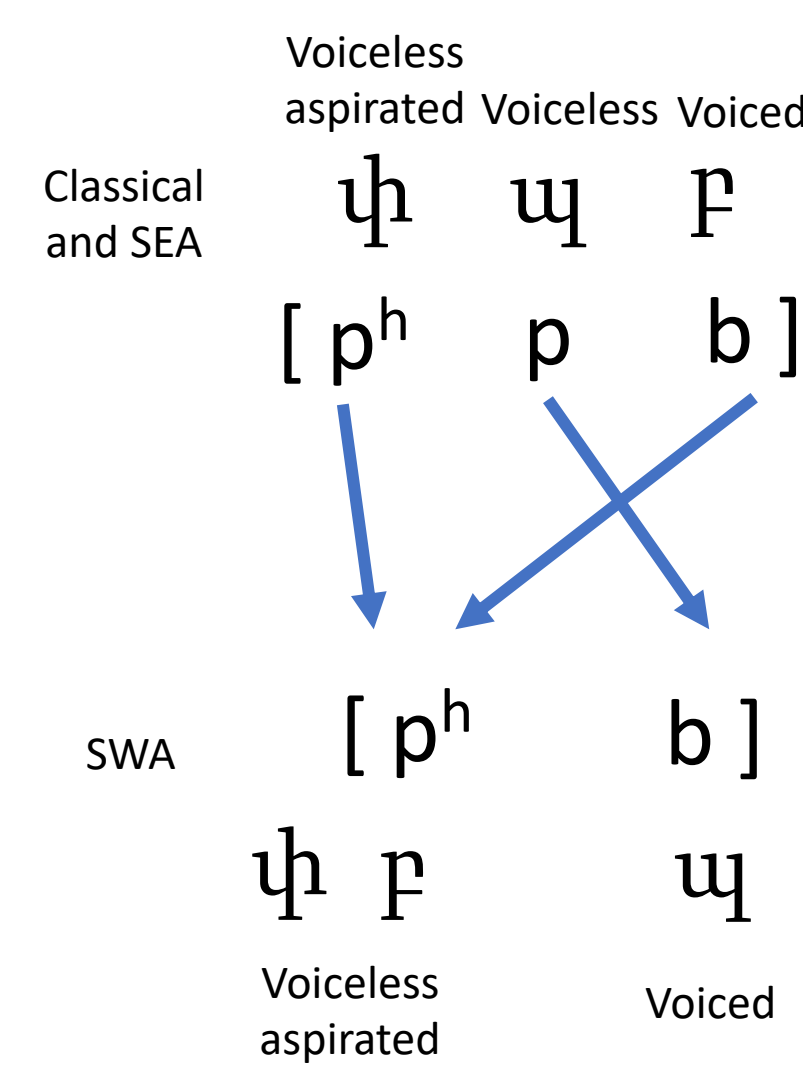
- Traditionally: Caucasus & Iran
- Currently: Republic of Armenia; Iran, post-Soviet countries, Europe and North America
- 36 phonemes : 30 consonants, 6 vowels
- **Three-way voicing system**
- Voiced; voiceless; voiceless aspirated

	Bil.	L-d.	Alv.	P.-Al.	P	Vel.	Uv.	G
Nas.		m		n				
Plos.		p ^h p b		t ^h t d				k ^h k g
Affr.				ts ^h ts dz				tʃ ^h tʃ dʒ
Fric.		f v		s z		ʃ ʒ		χ ʁ h
Appr.				l			j	
Tap				r				
Trill								

Armenian's orthography (in both variants) maintains a representation of three graphemes for each of these voicing sequences, making rule-based speech synthesis of either pronunciation feasible from the same text. However, producing text from speech poses a challenge when some acoustically identical inputs correspond to different graphemes in the two variants.

SEA	SWA
Voiced	Voiceless aspirated
Voiceless	Voiced
Voiceless aspirated	Voiceless aspirated

item	SEA	SWA	translation
< բառ >	[bar]	[p ^h ar]	'word'
< պար >	[par]	[bar]	'dance'
< փայտ >	[p ^h ar]	[p ^h ar]	'placenta'



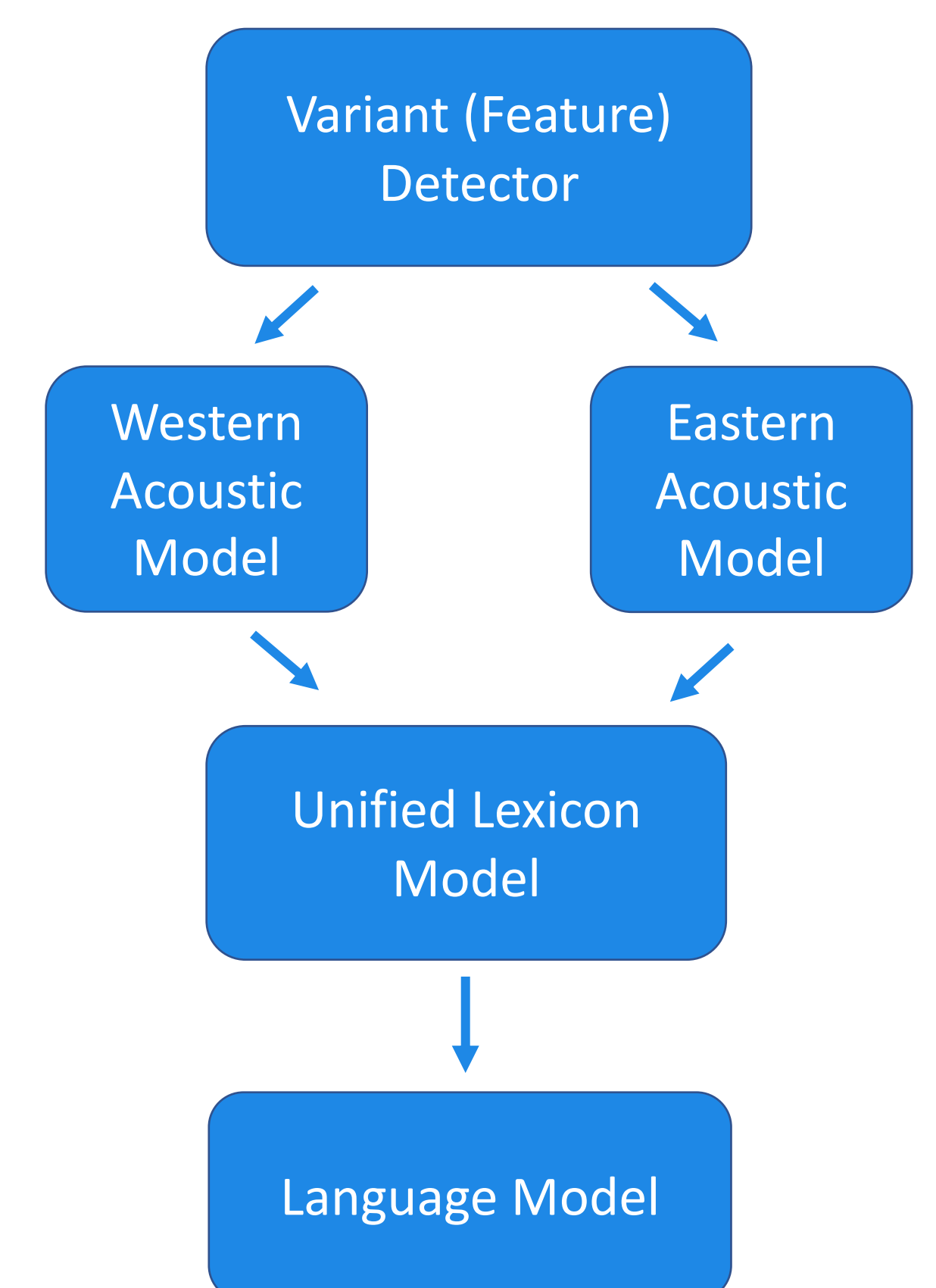
SOLUTIONS

To tackle this problem, we will first construct a traditional hybrid ASR model, with the eventual goal of comparing the efficacy of the system to an End-to-end (E2E) model. In these tests, we will use the manually transcribed and aligned data from the EANC and Rerooted corpora to determine the level of accuracy yielded by the two methods.

Below, we propose two pathways for processing audio data from the two dialects in one (hybrid) system.

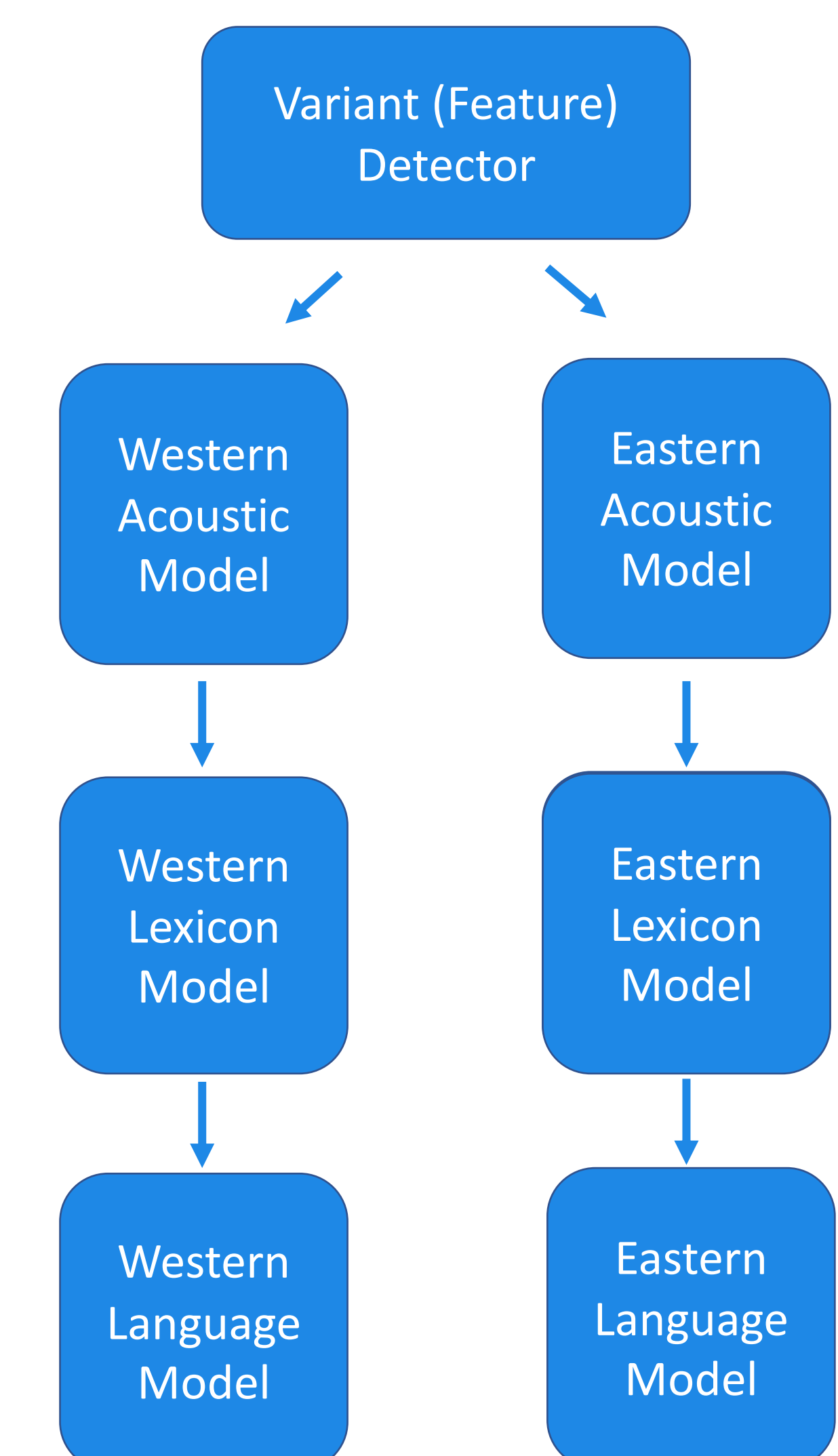
OPTION 1

Two separate acoustic models, which feed into a common lexicon model.



OPTION 2

Two completely independent paths for the two variants.



THE DATA

Our work to construct an ASR model for Armenian is conducted in the framework of the ANR-funded DALiH (Digitizing Armenian Linguistic Heritage) project, within which we expect to take advantage of two major transcribed audio corpora.

Rerooted Archives

REROOTED

- Western Armenian (speakers from Syria)
- corpus of interviews, monologues
- 120 hours of audio data
- 82 hours transcribed, and aligned by sentence

Eastern Armenian National Corpus (EANC)

EANC
Eastern Armenian National Corpus

- Eastern Armenian
- corpus of spontaneous, public and task-oriented discourse; online communications
- 774 total transcribed (audio) documents
- 3.5 million tokens

Samuel CHAKMAKJIAN
Սամուէլ Չաքմաքճեան
Doctoral student, Inalco SeDyL & ERTIM