

<p>Griots Speech Dataset</p>	<div style="text-align: right;">  </div> <p>The Griots Dataset is a combined Bambara audio, and semi-professionally transcribed, and translated into French speech dataset. This dataset was developed with the intention to build a Bambara ASR model.</p>	
<p>DATASET LINK</p>	<p>DATA CARD AUTHOR(S)</p>	
<p>https://github.com/robots-mali-ai/jeli-asr</p>	<p>Sebastien Diarra, RobotsMali: Owner</p>	
<p>https://zenodo.org/record/7296317</p>		
<p>Authorship</p>		
<p>Publishers</p>		
<p>PUBLISHING ORGANIZATION</p>	<p>INDUSTRY TYPE(S)</p>	<p>CONTACT DETAIL(S)</p>
<p>RobotsMali</p>	<p>Not-for-profit - Academic - Tech</p>	<p>Publishing POC: Michael Leventhal Affiliation: RobotsMali Contact: mleventhal@robotsmali.org Website: https://robotsmali.org</p>
<p>Dataset Owners</p>		
<p>TEAM(S)</p>	<p>CONTACT DETAILS</p>	<p>AUTHOR(S)</p>
	<p>Dataset Owner(s): RobotsMali Affiliation: RobotsMali</p>	<p>Michael Leventhal, PI</p>

	Contact: research@robotsmali.org Website: https://ai.robotsmali.org	Sebastien Diarra, Manager Mouktar Traore, Supervisor Alou Dembele, Supervisor														
Funding Sources																
INSTITUTION(S)	FUNDING OR GRANT SUMMARY(IES)															
GOOGLE, LLC	Grant provided to develop a functioning ASR model in Bambara; Also to improve the performance of an existing Bambara-French Translation Model.															
Dataset Overview																
DATA SUBJECT(S)	DATASET SNAPSHOT	CONTENT DESCRIPTION														
Cultural Data Social Commentary History Narrations	<table border="1"> <tr> <td>Size</td> <td>16 GB</td> </tr> <tr> <td>Length</td> <td>30 hours</td> </tr> <tr> <td>Utterances (Clips)</td> <td>29800</td> </tr> <tr> <td>Ave. Clips Length</td> <td>3.02 s</td> </tr> <tr> <td>Tokens</td> <td>300923</td> </tr> <tr> <td>Types</td> <td>62753</td> </tr> <tr> <td>M/F Speaker Ratio</td> <td>23/7</td> </tr> </table>	Size	16 GB	Length	30 hours	Utterances (Clips)	29800	Ave. Clips Length	3.02 s	Tokens	300923	Types	62753	M/F Speaker Ratio	23/7	The dataset consists of 1 hour recording of 30 griots. Each recording is focused on the historical and cultural records of Mali.
Size	16 GB															
Length	30 hours															
Utterances (Clips)	29800															
Ave. Clips Length	3.02 s															
Tokens	300923															
Types	62753															
M/F Speaker Ratio	23/7															
Dataset Version and Maintenance																
MAINTENANCE STATUS	VERSION DETAILS	MAINTENANCE PLAN														
Actively Maintained	Current Version: 1.0.0 Last Updated: 09/2022 Release Date: 08/2022	The repository is adjusted accordingly, to augment the data quality. Additional data in the domain is sought after. Versioning: Major.Minor.Patch.														

		<p>Semantic Versioning specification is adopted.</p> <p>Updates: This will occur on a cyclical basis, as we expect to increase the dataset with related works.</p> <p>Errors: Errors are addressed right when they are flagged. https://github.com/robotsmali-ai/jeli-asr/issues</p> <p>Feedback: Contact the manager, send us an email at research@robotsmali.org.</p>
--	--	--

Example of Data Points

PRIMARY DATA MODALITY	SAMPLING OF DATA POINTS	DATA FIELDS
Text Data Audio Data	https://github.com/RobotsMali-AI/jeli-asr/blob/master/README.md	<p>Utterance start (milliseconds)</p> <p>Utterance end (milliseconds)</p> <p>Utterance transcription</p> <p>Utterance translation</p> <p>Second elapsed</p> <p>Recording Fragment ID</p>

TYPICAL DATA POINT

[4640, 6460, "Alhamdulillah Rabil Alamiina", "Alhamduolilaye Rabil Alamina!", 1.82, "griots_r1_1"]

Motivations & Intentions		
PURPOSE(S)	DOMAIN(S) OF APPLICATION	Intended And/Or Suitable Use Case(s)
Safe for research use Safe for production use Training Inferencing	Machine Learning, Automatic Speech Recognition, Machine Translation	- Bringing the bambaraphone communities online, through cultural data.

		- Resource generation for a low-resource language.
Provenance		
Collection		
METHOD (S) USED	METHODOLOGY DETAIL(S)	SOURCE DESCRIPTION (S)
Crowdsourced - Paid	<p>Audio Recording Primary modality of collected data: Audio Data</p> <p>Transcription/Translation Platform: ELAN Primary modality of collected data: Text Data</p> <p>Dates of collection: 05 2022 - 08 2022</p>	<p>Audio: Griots. Oral stories recorded with professional recording equipment.</p> <p>Text: Linguists / Transcribers, using ELAN</p>
Use in ML or AI Systems		
Dataset Use (s)	Notable Features	Usage Example
<p>Training</p> <p>Testing</p> <p>Validation</p> <p>Fine Tuning</p>	<p>Tokens</p> <p>Types</p> <p>Utterances</p>	<p>https://github.com/RobotsMali-AI/jeli-asr/tree/master/asr</p>
Annotations & Labeling		
ANNOTATION WORKFORCE TYPE		
<p>Human Annotations (Non-Expert)</p> <p>Human Annotations (Contractors)</p> <p>Human Annotations (Crowdsourcing)</p>		

Human Annotators		
	ANNOTATOR DESCRIPTION	ANNOTATOR TASK
	<p>Transcription / Translation Use ELAN to transcribe audio fragments.</p> <p>Number: 20 persons team</p> <p>Level: Amateur linguist</p> <p>Editors Edited transcriptions and translations.</p> <p>Number: 5 persons team</p> <p>Level: Experts</p> <p>Language Distribution of annotators:</p> <ul style="list-style-type: none"> - French: 3 - Bambara: 2 	<p>Transcription / Translation: Transcribing audio fragments using ELAN and translating.</p> <p>Editors Take the ELAN file from the transcription step</p>
LICENSING		
LICENSE TYPE(S)	LICENSE BREAKDOWN	LICENSE PERMISSIONS
CC-BY-SA 4.0	<p>This work is licensed under the Creative Commons Attribution 4.0 International License.</p> <p>CC-BY-SA 4.0</p>	<ul style="list-style-type: none"> ● Share — copy and redistribute the material in any medium or format ● Adapt — remix, transform, and build upon the material for any purpose, even commercially. ● Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any

		reasonable manner, but not in any way that suggests the licensor endorses you or your use.
--	--	--

RobotsMali

Centre National Collaboratif pour l'Éducation en Robotique et en Intelligence Artificielle
Ecole Normale d'Enseignement Technique et Professionnel (ENETP)
Cité-Universitaire de Kabala
Bamako, Mali

Preliminary