

# Quantifying the Ethical Dilemma of Using Culturally Toxic Training Data in AI Tools for Indigenous Languages



**Pedro  
Domingues**



**Claudio  
Pinhanez**



**Paulo  
Cavalin**



**Julio  
Nogima**

*pcavalin@br.ibm.com*  
*csantosp@br.ibm.com*

IBM  
Research

# PROINDL

## AI technologies to strengthen Indigenous languages in Brazil



Claudio  
Pinhanez



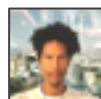
Paulo  
Cavalin



Luciana  
Storto



Thomas  
Finbow



Alex  
Cobbinah



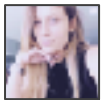
Julio  
Nogima



Isabel  
Gonçalves



Sarajane  
Peres



Nicole  
Dalmiglio



Majoi  
Gangora

IBM Research, Brazil

University of São Paulo



### Multidisciplinary Team

- 2 research scientists, 1 software engineer and 1 doctoral intern from IBM Research, Brazil
- 3 professors from the USP Dept. of Linguistics (Indigenous languages)
- 1 prof. from USP IT Dept. (robotics)
- 1 post-doc (USP Anthropology)
- 1 technical support staff
- 6 undergrad scholarship interns

### Student Collaborators



MISTI - MIT-Brazil



Insper (São Paulo)

<https://c4ai.inova.usp.br/research-activities-in-the-c4ai/#>



## focus areas

Development of **writing assistants** for text production and social media use based on Indigenous Language Models.

Desenvolvimento de **aplicativos WhatsApp e Android** para fornecer melhor suporte e acesso a recursos às comunidades.

Developing **translators from Brazilian Indigenous languages** to Portuguese by fine-tuning ML translators.

## prototypes



writing assistant



WhatsApp app



Android app

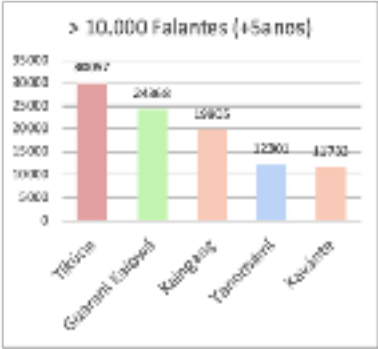
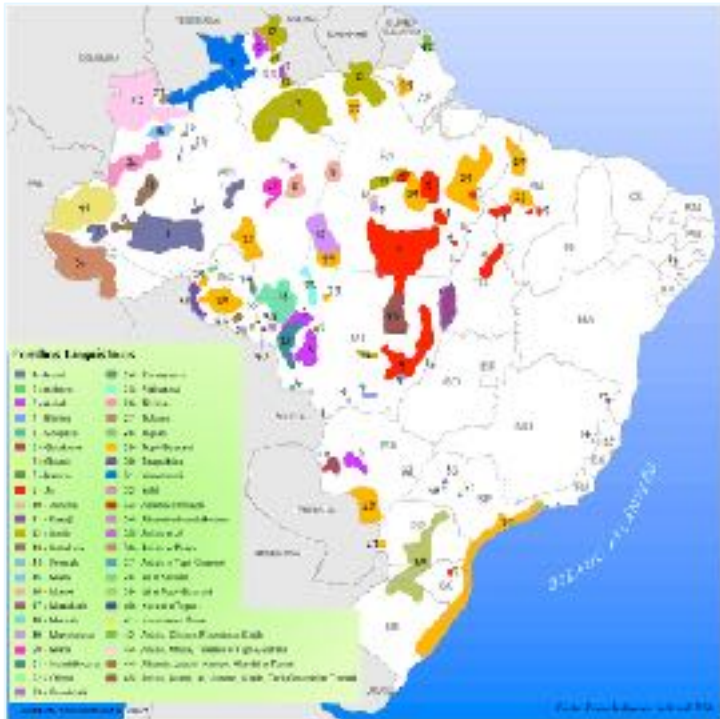


bidirectional translation of Indigenous languages

## community work



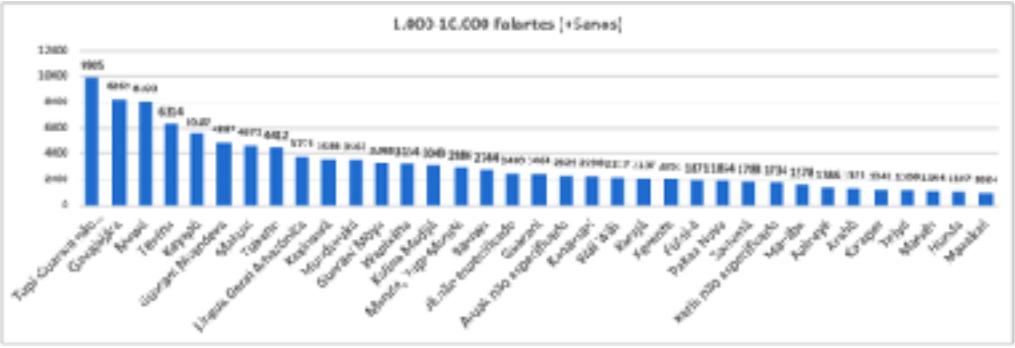
# The majority of the 200 Indigenous languages in Brazil are likely to disappear in the next 100 years



## Brazil 2010 census:

270 languages  
800K speakers, half in Ind. lands

>10K speakers: 5 languages  
1K< speakers < 10K: 35 languages



# The Bibles dataset: translations of the *New Testament* in 39 Brazilian Indigenous languages

Mathew 1:6. “and Jesse the father of David the king. And David was the father of Solomon by the wife of Uriah.”

**Apalai** => Jese mokyro tuisa Tawi zumy. Tawi mokyro Saromão zumy. Saromão eny mokyro Uria pytyÿpny.  
**Aprinã** => Xesee, Tawil iri, Tawil Iseakini asite Itxua. Tawil, Saromão Iri, Saromão Inoro mitxi Oriá Itanoro. Eereka  
**Asheninka** => ikanta Isaf tzmanake itomi Irinori, itakeri David. Iri pinkatharitaintsiri, ikanta pinkathari David rowaly  
**Culina** => Dsesie-Isahipa Dabicca abi. Najari Dabioa jodiidenicca tamine tojajari. Póhuapa Saromocca abi. Saromo  
**Desano** => Isal Davi pagw árllumí. Davipe Israe majarã tauro opu árllumí. Ige Uriá walçegu mewe marapore Ige m  
**Guajajara** => Zxe umuzãg tuwihaw Tawi a'e. Tawi umuzãg Xárumãw a'e. (Ihy Uri hemireko kwer maro hekon a'e.)  
**Guarani\_Mbya** => Jessé ra'y ma huvixawe Davi, ha'e huvixa Davi ra'y ma Salomão, Urias ra'yrykue pi'a va'e.  
**Jamamadí** => ãfã batimatamona amaka Sese. Saomao batimatamona amaka Tañ. Iaromao matã Oriã fatã tohibehã  
**Kadiwéu** => Jessé jigjãa eliodi Davi, Davi anijo me inionigi-eliodi. Davi eliodi Salomão, eliodi Salomão nagejo anijo  
**Kaingang** => Kÿ Jessé vÿ pã'i mág tÿ Davi ên han. tÿ pã'i mág tÿ Davi vÿ Salomão ha, tã prũ tÿ tÿ Urias prũ ja ên fi ki  
**Karajá** => Davi háwyy heka Uriá háwyy juhúu raremyhÿ. Iribi Davi-wana roire. Salomã Iróre rare. Tii heka Roaboã  
**Macushi** => Maropai Jessé wani'pi Davi yun pe, pita esa' pe tñw'sen.  
**Maxakali** => hu Namix mûg tak. Yã hömã Namix ta 'hyaet yög tikim'ün xohix xat. Ha Namix te 'Ovit xetut mûg, ha X  
**Nadëb** => Jessé t'aah kã, Dawi, êr wahé makú sa wahé m'aa paah. Dawi t'aah, Saromãw, Saromãw ÿÿn Urija hád nãn  
**Nambikuara** => Je'sah'ta'ki'ha'ty' Ta'vi'yah'lo'au' ta'hoai'he'ta'. A'noe'jah'tai'na'ti'sa'e'p' nõ'kg'ho'a' a'wo'ka'e'n'y  
 a'foe'ho'ka'ho'ka' a'wã'ho'ka'ho'ka' Uri'ho'ka'ho'ka' a'ho'ka'ho'ka'ho'ka' a'ho'ka'ho'ka'ho'ka'. A'noe'ho'ka'ho'ka' a'ho'ka'ho'ka'ho'ka' S.  
**Palikúr** => Igme amekene Jessé gikamkayh ukiparawiy amekene Davi. Igme amekene Davi xuwehe amekene Uriyas  
**Parecis** => Jessé atyo Davi kaisani hoika Davi atyo ekohaseti kalorexi tyona. Hatyoseta Davi atyo Salomão kaisan  
**Rikbaktsa** => Takino Dawizo. Iwa taparakta Abarão tsekokatsa niaha. Iwaze Saromae ta Dawi tse ije tapara Urias ok  
**Sateré** => Mi'i kawyi IESE imoherep alporekuat s'awy'iwuanuat TAWI. Mi'i hawyi norekuat koro TAWI kaipyi tuwe  
**Tenharim** => Jesséva'ea Daviva'ea po'ria. Daviva'ea israelitasva'ea nuvihavuhuva'ea. Daviva'ea Salomõova'ea po'ria.  
**Terêna** => Enejoneko Njese, énomone itúko há'aneko Ndávi, náti mekúke.  
**Ticuna (Peru)** => Rú nûma ga Ichaxi rû ãÿÿgarú gi Dabinatú nibí. Rú nûma ga ãÿÿgarú ga Dabi rû Charomõúnatú nã  
**Tucano** => Isal' udio masã wlogã Davi pacu nice riwã. Davi Urias nemo ni'co me'rã Salomõrê põ'rãtice riwã.  
**Wapishana (Guiana)** => Jesse dani uo King David. Aizii di'i Jesus Christ dokozu-daurnao da'ana'o King David di'ki'o  
 [Solomon dano Uriah daiani-daun];

Name	Indigenous Languages				# Aligned Sentences		
	Acron	Branch	Family	Speakers	Train	Test	Total
Bororo	bor	Macro-Jê	Bororo	1275	1801	200	2000
Apinãwê	apn	Macro-Jê	It	827	827	76	903
Calungsã	lgs	Macro-Jê	It	1900	9695	913	10513
Ejagob	tsu	Macro-Jê	It	5500	2629	512	3141
Enawtse	xau	Macro-Jê	It	1200	1275	340	1615
Karajã	krj	Macro-Jê	Enaja	1819	2808	587	3615
Matsigenka	mkg	Macro-Jê	Matsigenka	1248	5588	586	6174
Rikbaktsa	rbs	Macro-Jê	Rikbaktsa	10	1589	712	2301
Miwã	mwa	Basí	Miwã	800	638	97	735
Mundurukã	mva	Basí	Mundurukã	1363	2727	126	2853
Guajajara	gja	Basí	Tupi-Guarani	8269	4996	934	5890
Guarani (West-Brazil)	gwr	Basí	Tupi-Guarani	NA	5203	97	5270
Guarani (East-Brazil)	geu	Basí	Tupi-Guarani	NA	5269	90	5359
Guarani (Central)	gca	Basí	Tupi-Guarani	34908	9694	479	35381
Guarani (Mbya)	gmb	Basí	Tupi-Guarani	1248	6349	97	7114
Guarani (Paraguay)	gpa	Basí	Tupi-Guarani	NA	5146	97	5243
E'wãpã	ewp	Basí	Tupi-Guarani	1261	3388	437	3816
Katãmbã	kat	Basí	Tupi-Guarani	128	2387	286	2613
Mhẽngatu (Gua)	gma	Basí	Tupi-Guarani	3771	5685	690	9746
Tenharim	ten	Basí	Tupi-Guarani	73	3395	644	4072
Iwãwãwã (Guarani)	iwa	no branch	Arãwã	217	4799	715	5271
Cãtãna Modjji	mod	no branch	Arãwã	3543	4339	690	8572
Pãwãwã	pwv	no branch	Arãwã	160	2673	379	3252
Apãrã	apv	no branch	Arãwã	54	6329	97	6376
Pãtãrã	ptv	no branch	Arãwã	424	6147	968	7139
Pãrãrã	prv	no branch	Arãwã	52	6381	97	6478
Tãrãrã	trv	no branch	Arãwã	6314	6381	97	7092
Wãgãrã	wgv	no branch	Arãwã	194	568	86	648
Enãwãrã	enw	no branch	Guãrãrã	249	4523	790	5262
Apãrã	apv	no branch	Karã	251	5548	97	6196
Bãrãrã	brr	no branch	Karã	173	4809	317	5126
Hõrãrãrã	hrr	no branch	Karã	53	4249	67	4316
Uãkãrã	ukv	no branch	Karã	429	490	94	583
Mãrãrã	mrv	no branch	Mãrã	58	5232	813	5373
Nãrãrãrãrã	nrv	no branch	Nãrãrãrãrã	201	274	84	313
Enãwãrãrã (Perã)	enw	no branch	Para-Tucano	1368	3146	187	3661
Tãrãrã	trv	no branch	Tucano	4412	279	848	4299
Yãwãrãrãrã	yrv	no branch	Yãwãrãrãrã	1200	1785	198	2483
Tãrãrã	trv	no branch	no family	2827	2827	282	3483
<b>TOTAL</b>	<b>29</b>	<b>3</b>	<b>95</b>	<b>109001</b>	<b>150025</b>	<b>2580</b>	<b>108603</b>

# Using *The Bible* in Indigenous contexts: Issues and Concerns

- *The Bible* is a religious text and sacred to many people in the world, including Indigenous people, and therefore should be treated with great respect and care.
- *The Bible* is also connected to negative aspects of past and present colonial history of Indigenous peoples, in particular to the effort to convert them to Western religions, in particular to Christianity.
- The translations have established *orthographies of domination* into many Indigenous language. There also many issues about the quality of the translations and common problems of “Europefication” of languages [Franchetto, 2008].

[Bruna Franchetto. 2008. The war of the alphabets: indigenous peoples between the oral and the written. *Mana*, 14(SE):31–59.]

***“the Bible was a tool for the colonization process [...working] hand-in-hand in the exploitation, subjugation, and continued oppression of the Indigenous Peoples of the U.S.”***

[Chris Mato Nunpa. 2020. The great evil: Christianity, the bible, and the Native American genocide.]

***“at the beginning of the colonization process two tools of genocide were forced upon Native people: the bottle and the bible.”***

[Stormy Ogden. 2005. The prison-industrial complex in Indigenous California. In *Global lockdown: Race, gender, and the prison-industrial complex.*]

# Ethical issues with Indigenous data



Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence (IJCAI 23)  
Specialized track for Covid

## Balancing Social Impact, Opportunities, and Ethical Constraints of Using AI in the Documentation and Vitalization of Indigenous Languages

Cassio S. Pinheiro, Paulo Cavalari, Maria Vasconcelos and Julio Magina  
IBM Research, Brazil

[cassio.sp, paulo.cavalari, mariava, magina]@br.ibm.com

### Abstract

In this paper we discuss how AI can contribute to support the documentation and vitalization of Indigenous languages and how that involves i) delicate balancing of ensuring social impacts, exploring societal opportunities, and dealing with ethical constraints. We start by surveying previous work on using AI and NLP to support critical activities on strengthening indigenous and endangered languages and discussing key limitations of current technologies. After reviewing basic ethical constraints of working with Indigenous languages and communities, we propose that creating and deploying language technology efficacy with and for Indigenous communities leverages AI researchers and engineers to address some of the main shortcomings and criticisms of current technologies. Those ideas are also explored in the discussion of a real case of development of large language models for Brazilian Indigenous languages.

### 1 Introduction

Tradition and progress are often in conflict in Indigenous communities and one of its most common battles is in strengthening the use of their own languages. We argue in this paper that using Artificial Intelligence (AI) and, particularly, Natural Language Processing (NLP) technologies to support

social impacts of AI is one of the central themes of our paper and hereby we assume it to mean using AI technology to contribute to the solution of sociolinguistic-economic problems of underserved and vulnerable communities, according to the needs expressed by them, respecting their social and cultural context, and, whenever possible, in projects led by them.

In fact, creating and deploying technology to be used in Indigenous communities must follow ethical guidelines, as discussed in Kitching et al., 2017; Strick et al., 2015). These are in stark contrast with traditional practices of AI, such as reliance on big data, data extractivism, and colonial thinking (Strickson, 2021), ideas as epitomized by the motto “Technology as colonization” adopted by the Indigenous communities as a guideline for any language initiative with them; any work, even in research projects, must be done with the community and for its benefit and in a sustainable manner.

We start this paper recognizing the importance of indigenous peoples and cultures in the world context. We follow with some data and definitions about endangered languages, a discussion about the value of language diversity, and with an overview of benefits and challenges of documentation and vitalization of indigenous languages. We then revisit some of the limitations of current NLP technologies to deal with Indigenous languages including issues with large language models (LLMs) such as BERT and GPT-3.

After discussing ethical issues and guidelines when working with Indigenous communities, we examine a research initiative conducted by some of the authors of this paper to use

IJCAI 2023

Can **culturally toxic data** increase the performance of translation models made for extremely low-resource languages?

- Given that this type of data can raise ethical concerns depending on the context in which it is applied, its use can be beneficial?



# Datasets

## Dictionary dataset:

- 1,022 short stories aligned sentences in Guarani Mbya (gun) and English (eng);
- 245 sentences from pedagogical material;
- 2,230 sentences from Dooley's Lexical Guarani Mbya dictionary.
- The last two sources were aligned with Portuguese (por), so we used Watson to translate from por to eng.
- After a data cleaning, unicodes normalization, the Dictionary dataset was splitted into 3,155 and 300 gun-eng aligned sentences for training and test.

## Bibles dataset:

- 39 translations of the Bible's New Testament, totaling 188,033 BILs-eng aligned sentences.

Divided into 3 training datasets:

- **Bilingual** (only Guarani Mbya): 6,340 training pairs;
- **Tupi Guarani Family** (10 BILs): 43,869 training pairs;
- **All BILs available** (39 BILs): 162,225 training pairs.

For testing, the Matthew chapter from Guarani Mbya New Testament (970 aligned sentences).

# Models

zeroshot: WMT19 model;

mbya: WMT19 model finetuned with bilingual data;

TGf: WMT19 model finetuned with Tupi Guarani family data;

all: WMT19 model finetuned with all BILs available data;

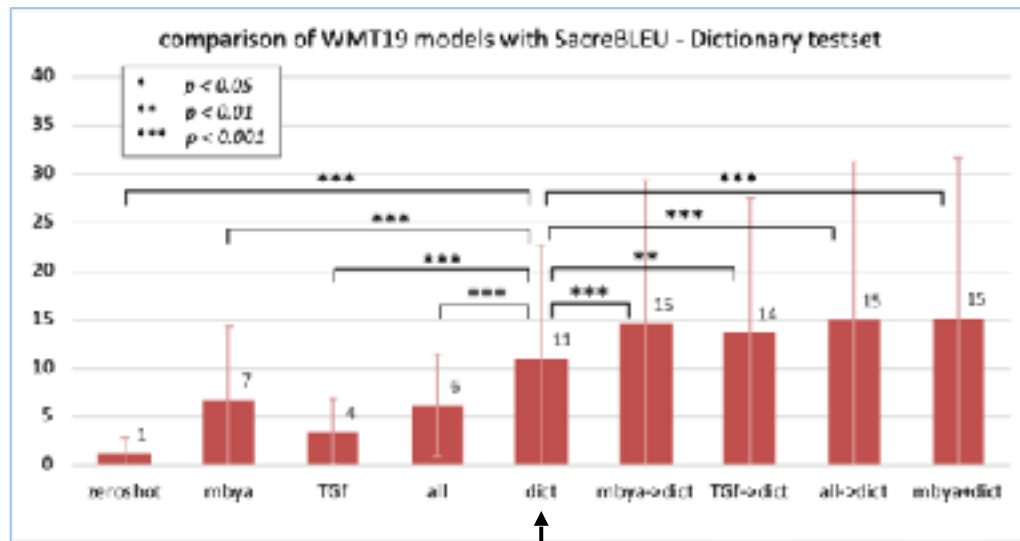
dict: WMT19 model finetuned with dictionary dataset;

mbya->dict: mbya model finetuned with dictionary dataset;

TGf->dict: TGf model finetuned with dictionary dataset;

all->dict: all model finetuned with dictionary dataset;

mbya+dict: WMT19 model finetuned with a merge of bilingual data and dictionary dataset.



no Bible data

# Qualitative analysis

models	mbya				dict				mbya+dict			
	EXPECTED OUTPUT FROM TEST SET	mbya OUTPUT	Sacre BLEU	BLEU	chrF	dict OUTPUT	Sacre BLEU	BLEU	chrF	mbya+dict OUTPUT	Sacre BLEU	BLEU
long ago there lived a giant	there were two men in the crowd.	83	43	30	long ago there lived a giant.	100	100	100	long ago there lived a giant	100	100	100
one day he went to the woods again.	then one of them went to the other side of the lake.	83	35	32	one day he went again to the woods.	46	91	70	and so one day he went again to the woods.	36	75	73
when he arrived at his house, he said to his wife, "can there be anyone who can hunt like me?"	when he came to the tomb, he said to his mother, how can I not know where I am.	17	47	30	when he got home, he said to his wife, "could it be that I'm from here?"	37	46	30	when he got home, he said to his wife, "could I find the ring?"	32	42	59
when he fell, he hit his back on the ground and died then and there.	so then, how much more will the earth bear down on him than the earth will bear down on him.	3	23	21	as he fell, he hit his forehead on the rock.	37	40	30	and as soon as he touched the ground, he died too.	7	41	50
years ago when I was a child, I didn't know the language of non-Indians.	I have not been able to speak the word of the one who sent me into heaven.	3	18	15	years ago when I was a lot younger, I didn't know what to do with the books.	25	55	40	years ago when I was a child, I did not understand the meaning of portuguese.	52	60	54
when my brother went, saw a snake.	when he came to my house, he saw me.	6	37	16	my brother went out to see the snake.	27	56	56	my brother went and saw the snake.	24	73	61
there comes an inhabitant of the bare village.	you are one of the twelve living creatures.	10	25	20	there comes the hare from the bare.	15	46	30	there comes the rabbit village.	15	46	63
he grabbed him by his arm	so he went up to heaven with his brother.	5	15	12	he took his brother to law there.	7	17	12	then he took half of the Indian in the sky.	4	6	9
when evening came, the birds were singing and singing, but the indian was still stuck.	but the spirit of the spirit is in the spirit, and the spirit is in the spirit.	5	25	17	and then it was the time to eat the birds, both of which were Indians.	4	41	35	and the one who drinks the spirit remains in it, though the spirit remains.	5	22	17
you changed religiously when you were to get me.	if I am a believer, it will be a believer in you.	3	5	9	if you guys believe me, it will believe you.	5	14	11	you will delude me even more.	6	11	12
who come with lower and higher people;	and all who are in the world and all who are in the world	4	16	17	have a lot of faith in him.	6	12	10	low-cost and high-cost carriers also must go.	7	22	25

examples of contaminated (in red) and non-contaminated outputs of the models mbya, dict, and mbya+dict

# All contaminated verses we could find....

expected output	generated by mbya-dict	BLEU	chrF
he killed three tapirs	he killed three of the <b>jesus</b> ,	24	61
dust.	a <b>bagel</b> of dust.	21	64
let's take out the stomach of this pig.	<b>i will pronounce</b> this pig.	18	24
then he saw something like a man open a door in the rock cliff.	then just as <b>the stone was coming out of the tomb</b> , something like a man opened the door.	14	49
the "claw-man" took the indian home to be her husband.	this man took the indian <b>and brought him to life</b> .	13	36
but his foot stuck too.	once again <b>he washed his feet</b> .	8	13
don't spill the tea or do you want to wet the bed completely?	do you not <b>untie the strap of your sandak</b> or sandak?	7	17
the pernilongs bit the one who was sleeping.	the <b>dove dove</b> .	5	6
will it be by chance that bad things happen to us?	<b>have we not turned a blind eye to evil?</b>	4	13
if he wasn't sick, it wouldn't have come.	<b>if i hadn 't been born, i would have never been born</b> .	4	23
if you had treated me i would have been cured.	<b>if i die, i die;</b>	4	4
when evening came, the birds were singing and singing, but the indian was still stuck.	<b>and the one who drinks the spirit remains in it, though the spirit remains</b> .	3	17
i said something like, "you what came already what laugh at guarani."	he was very pleased with the way he talked about it: "we could have bought a hat that would belong to <b>jesus</b> ."	3	17
sawing cable	small tree with adjective subordinate <b>prayer;</b>	0	11

14 of 300 (4.7%) presented some level of contamination (including 2 direct "Jesus" citations)

# Final discussion

The study demonstrates that using culturally toxic data **can significantly improve** the performance of LLM-based translators for ULR languages

- 30% improvement

The use of culturally toxic data can lead to potentially problematic outputs

- 4.7% contamination

**Careful consideration and communication with the communities involved**

We suggest its use only in **controlled situations** to mitigate negative effects, emphasizing the importance of community involvement and decision-making in the use of such tools.

The results highlight the need for more diverse training data, - future efforts to involve academic works, community-created data, and synthetic data generation in collaboration with linguists and language experts to **enhance the translation quality while respecting cultural sensitivities.**

# Thank You !



Paulo  
Cavalin



Claudio  
Pinhanez



Pedro  
Domingues



Julio  
Nogima

## Contributions

- Understanding that the Bible is toxic in Indigenous contexts.
- Quantification of the impact of the use of culturally toxic Bible data in the creation of Transformer-based Indigenous language models.

*Contact:*

*[pcavalin@br.ibm.com](mailto:pcavalin@br.ibm.com)*

*[csantosp@br.ibm.com](mailto:csantosp@br.ibm.com)*