



ansari4@purdue.edu; batra17@purdue.edu; lee3999@purdue.edu; chen3876@purdue.edu; gpahuja@purdue.edu; sharm536@purdue.edu; lanhamm@purdue.edu

## ABSTRACT

There is a critical need for optimizing the seed data needed for creating a widely acceptable machine translation model for low-level local languages. This research addresses the seed data concern by determining an optimized order of seed data which results in both more accurate and quicker translations as compared to a random order. This is achieved by dividing data from large translation project into various combination of test and train sets and achieve a BLEU score on the test data in the least amount of time and with the least number of iterations.

## INTRODUCTION

HIGH



\* Translation cost is a significant limiting factor on the pace and availability of translated important content.

HIGH

LARGE

- \* There are only handful of available methods to achieve the accuracy with minimum seed data usage.
- \* This project will demonstrate the 'high-accuracy low-data' dependent' algorithm that can be generalized and scaled across different languages to create effective translation for low-level local languages.

# **RESEARCH OBJECTIVES**

## **Our research focuses on answering the following questions:**

- What are the most important factors that play a role in optimizing the seed data for language translation?
- What kind of role does supplemental data pertaining to similar semantic domain play in optimization?
- How can business across world utilize this research to reduce translation cost?

# LITERATURE REVIEW

While all relevant studies focus on only machine translation or optimization (e.g. dynamic selection of seed data), none of them focused on optimizing seed data through ordering them differently based on semantic aspects.

Study	Ordering based on semantic meaning	Machine Translation	Optimization		
Our methodology (2022)	YES	YES	YES		
Liwei Wu, Shanbo Cheng, Mingxuan Wang, Lei Li (2021)		YES			
Xinyi Wang and Graham Neubig (2019	)	YES	YES		
Laura Martinus and Jade Abbott (2019	)	YES	YES		
Marline van der Mans Arianna Ricazza					
and Christof Monz (2017)	ι,	YES			
Ranathunga, S., Lee, E. S. A., Skendul M. P., Shekhar, R., Alam, M., & Kaur, R (2021)	İ,	YES	YES		
STEP 1 Create and clean parallel	corpora				
Book Chapter Verse Sequence	Engli	sh	Javanes		
0GEN11In th	e beginning God created the heavens and	t Ing jaman kawit	an Gusti Allah nitahake langit		
1 GEN 1 2 1 Now t	he earth was formless and empty, darkness	S Anadene bumi mau ca	ampur-bawur lan suwung, sega.		
29952 REV 22 20 1144 H	e who testifies to these things says, "Yes,	I Kang nglairake pas	seksen bab iki mau kabeh ngan		
29953 REV 22 21 1144 The g	race of the Lord Jesus be with God's peop	Sih-rahmate Gu	sti kita Yesus Kristus nunggila .		
STEP 2Get diversity & depth score for each chapter and split into train & test datasetDiversity & depthSplit $\square$ <					
			& depth		
Train 1 (diversity) Chapter 185 Chapter	181 Chapter 910 Cl	hapter 555 Cha	apter 776		
Train 2 (depth) Chapter 555 Chapter	124 Chapter 299 Cl	hapter 940 Cha	apter 24		
STEP 4 Run each subset of train	set into JoeyNMT mod	lel and get test	BLEU scores		
Train 1 (diversity)0-200 verses	JoeyNMT		BLEU score 1 BLEU score 2		

# **Intelligently Ordering Machine Translation Seed Data** to Improve Local Language Translation

Amaan Ansari, Devansh Batra, Jai Woo Lee, Paul Chen, Gagan Pahuja, Manideep Sharma, Daniel Whitenack, Matthew A. Lanham Purdue University, Krannert School of Management

tudy	Ordering based on semantic meaning	Machine Translation	Optimization
dology (2022)	YES	YES	YES
Shanbo Cheng, ang, Lei Li (2021)		YES	
raham Neubig (2019)		YES	YES
nd Jade Abbott (2019)		YES	YES
ees, Arianna Bisazza, of Monz (2017)		YES	
ee, E. S. A., Skenduli, , Alam, M., & Kaur, R. 2021)		YES	YES

Verse	Sequence	English	Javanese
1	1	In the beginning God created the heavens and t	Ing jaman kawitan Gusti Allah nitahake langit
2	1	Now the earth was formless and empty, darkness	Anadene bumi mau campur-bawur lan suwung, sega
20	1144	He who testifies to these things says, "Yes, I	Kang nglairake paseksen bab iki mau kabeh ngan
21	1144	The grace of the Lord Jesus be with God's peop	Sih-rahmate Gusti kita Yesus Kristus nunggila









# **STATISTICAL RESULTS**

### Test BLEU scores - div. vs seq.







## **BUSINESS IMPACT**