# Machine Translation and Transliteration
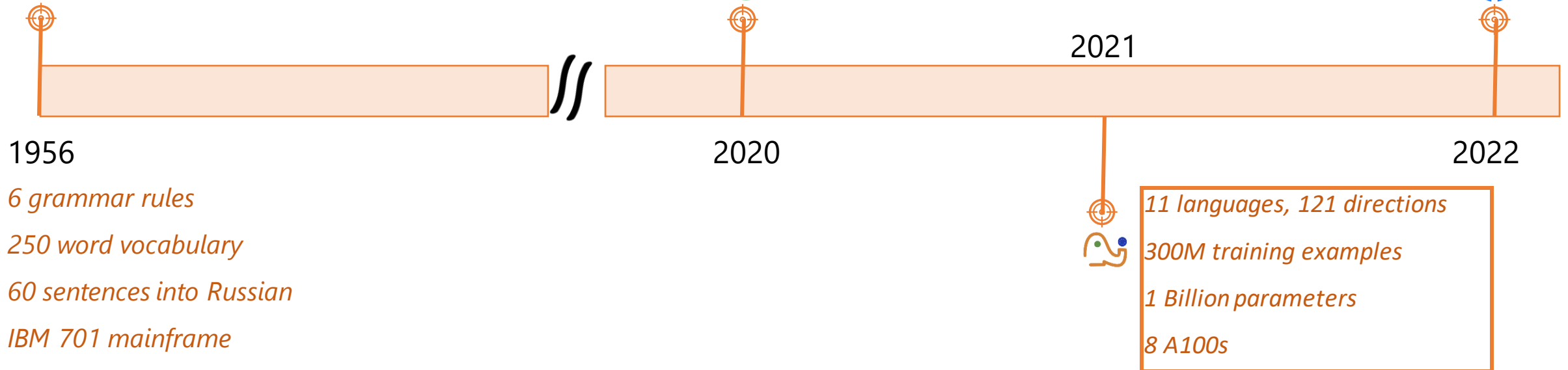
https://ai4bharat.iitm.ac.in/translation

https://ai4bharat.iitm.ac.in/transliteration

Workshop on 28th July 2022, IIT Madras

# Machine Translation: A brief history



Sequence to Sequence Learning with Neural Networks

Attention Is All You Need

No Language Left Behind: Scaling Human-Centered Machine Translation

The Mathematics of Statistical Machine Translation: Parameter Estimation

GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding

2015

2021

1956

1991

2014

2017

2020

2022

Neural Machine Translation by Jointly Learning to Align and Translate

*Samanantar*: The Largest Publicly Available Parallel Corpora Collection for 11 Indic Languages

# Machine Translation: A brief history

100 languages

13B training examples

1 Trillion parameters

2048 TPUs*

200 languages

XX training examples

YY parameters

ZZ A100s

2021

1956

2020

2022

6 grammar rules

250 word vocabulary

60 sentences into Russian

IBM 701 mainframe

11 languages, 121 directions

300M training examples

1 Billion parameters

8 A100s

# Machine Translation: A brief history

100 languages

13B training examples

1 Trillion parameters

2048 TPUs*

200 languages

XX training examples

YY parameters

ZZ A100s

2021

1956

2020

2022

6 grammar rules

250 word vocabulary

60 sentences into Russian

IBM 701 mainframe

11 languages, 121 directions

300M training examples

1 Billion parameters

8 A100s

# What is the recipe or modern NMT?

**DATA**

En　　हि

A large number of parallel sentences between En and Indic languages

**MODELS**

Large scale models with innovations specific to Indic languages

**EVALUATION**

Robust evaluation with diverse benchmarks and reliable evaluation metrics

**WEB SOURCES**

**Comparable**

**Non-Comparable**
**(Monolingual, billion tokens)**

**Machine Readable**

**Non-Machine Readable**

En हि

En हि

En हि गु ... ...

**Principle:** Curate data from the web, manual collection is too expensive and time consuming
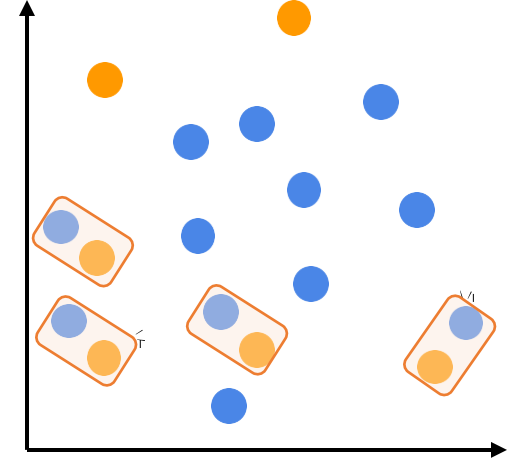
# How did we solve the data problem?



En

https://mykhel.com/
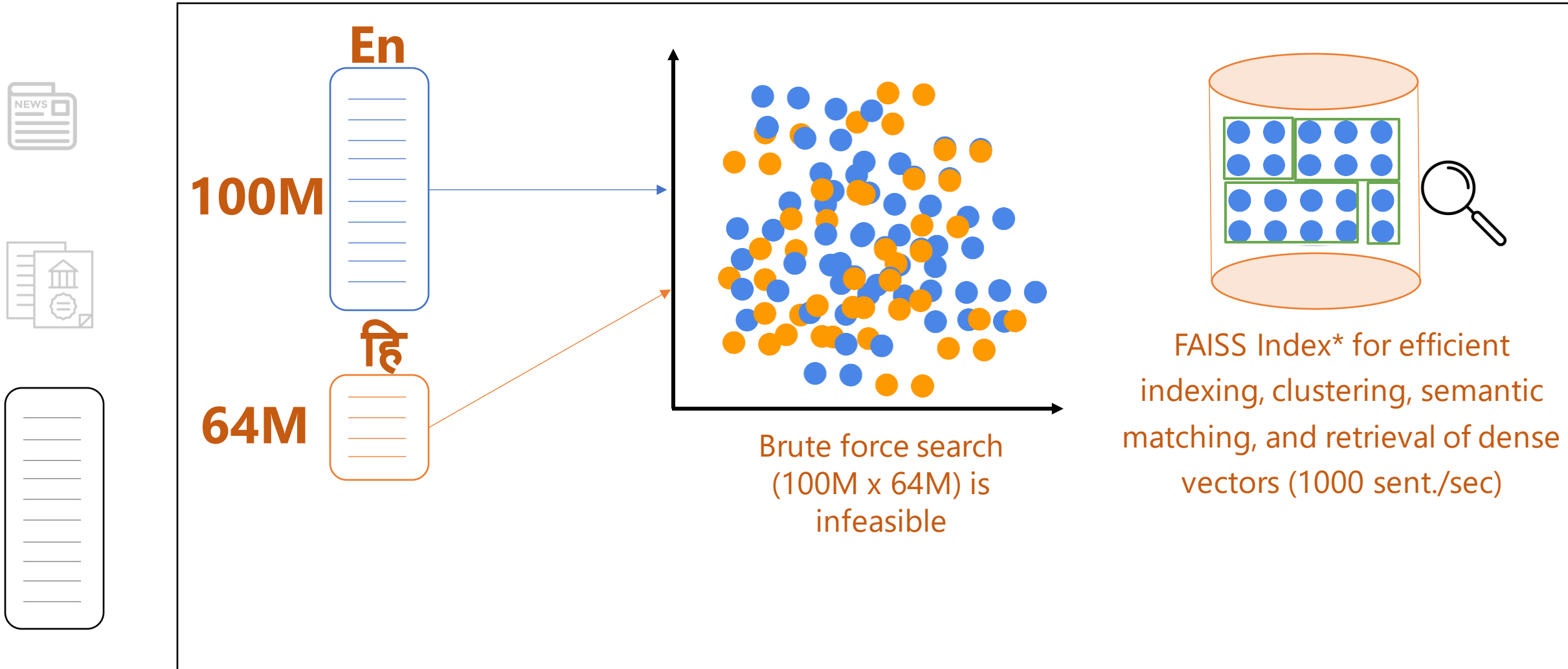
Jan 2020

ह

https://hindi.mykhel.com/

Shared multilingual space

24 such news sources considered in this work with data from 2010 onwards

**En**

**100M**

हि

**64M**

Brute force search
(100M x 64M) is
infeasible

FAISS Index* for efficient
indexing, clustering, semantic
matching, and retrieval of dense
vectors (1000 sent./sec)

# How much data did we collect?



**33M parallel sentences mined from the web (3X improvement)**

# What is the model that we use?



## Highlights

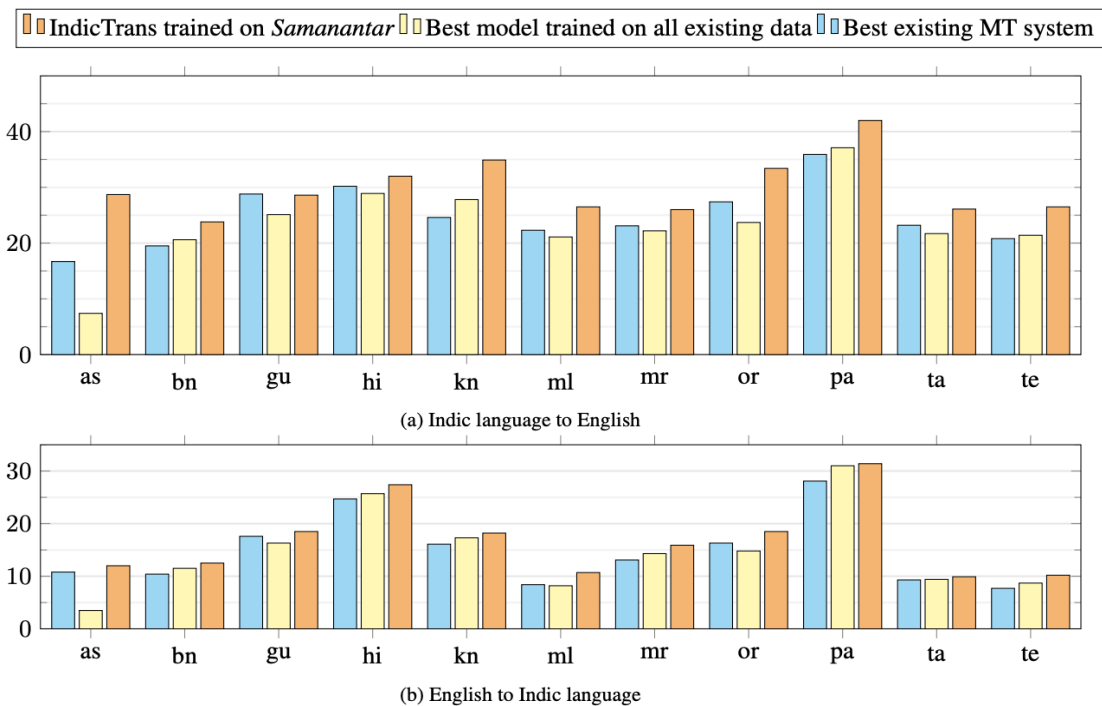Joint multilingual model for 11 Indic languages

Single script (enables transfer, reduces vocabulary)

3 models: En-X, X-En, X-X

6 encoder layer, 6 decoder layers, 16 heads/layer
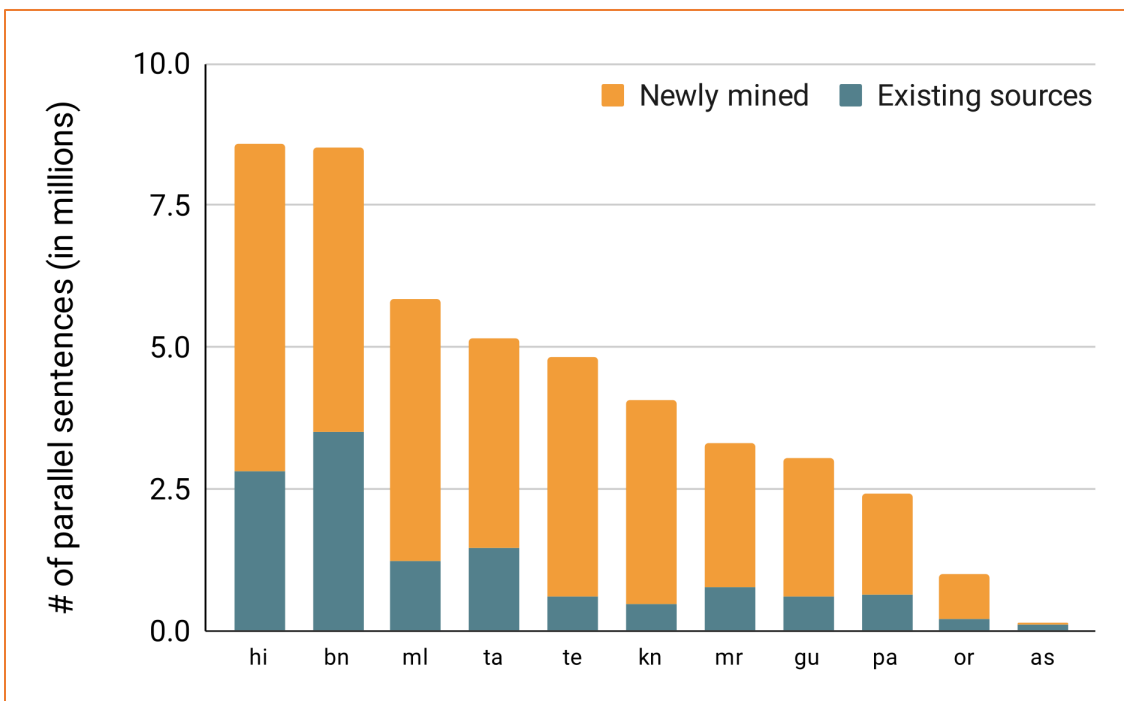
# What is the model that we use?



(a) Indic language to English

(b) English to Indic language

Legend: IndicTrans trained on *Samanantar* | Best model trained on all existing data | Best existing MT system
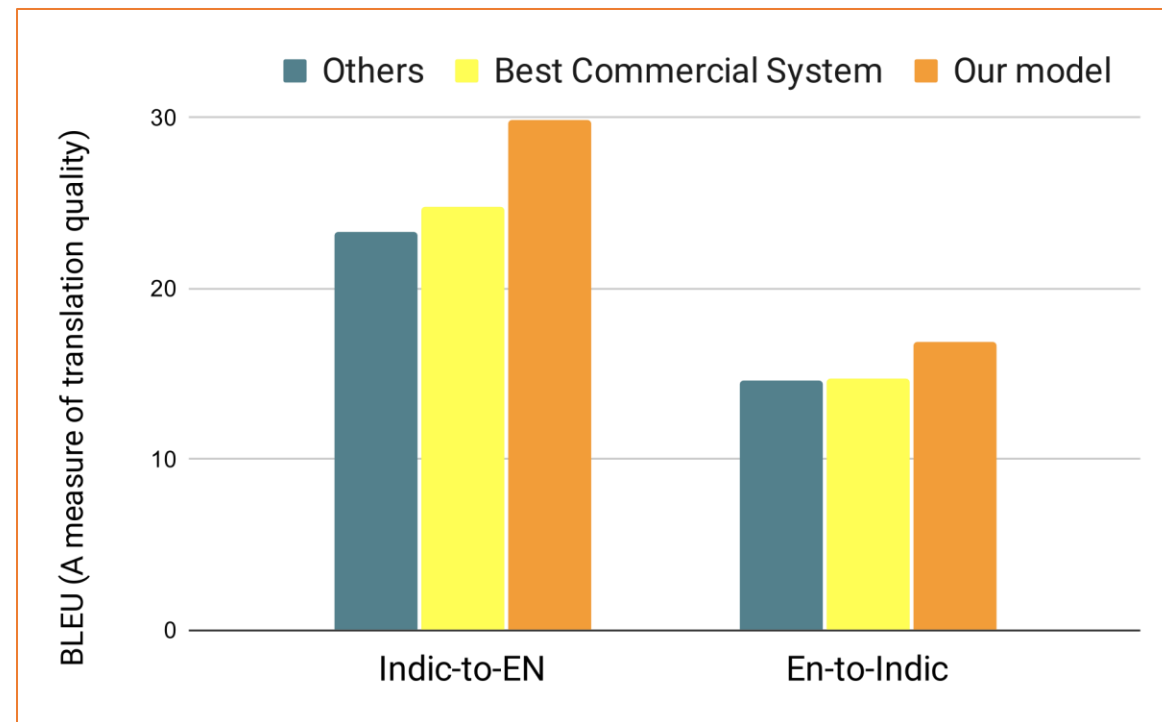
## Highlights

State of the art performance*

Gains are higher for low resource languages

# Summary



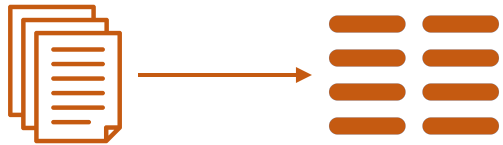**Samanantar**: Largest Parallel Corpus for
Indic Languages

https://ai4bharat.iitm.ac.in/samanantar

**IndicTrans:** State of the art translation
models for En-X and X-En

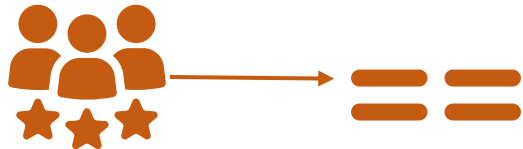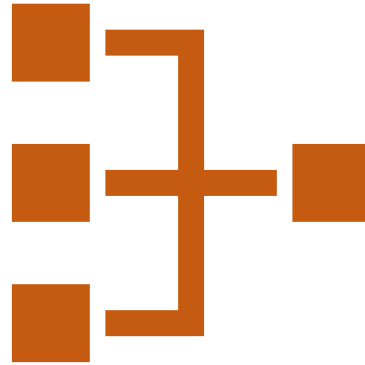https://ai4bharat.iitm.ac.in/indic-trans

# What next?

**DATA**

Mined data from diverse sources
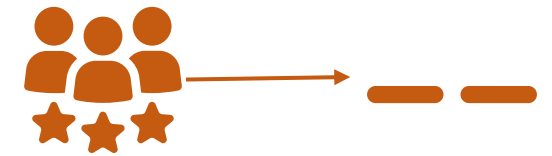
Manual training data

**MODELS**

Multilingual adapters

Simpler attention mechanisms for Indic-Indic translation

**EVALUATION**

Diverse Benchmark

Better evaluation metrics