# Revisiting Low Resource Status of Indian Languages in Machine Translation

Jerin Philip[*]
IIIT Hyderabad
Hyderabad, Telengana
jerin.philip@research.iiit.ac.in

Shashank Siripragada[*]
IIIT Hyderabad
Hyderabad, Telengana
shashank.siripragada@alumni.iiit.ac.in

Vinay P. Namboodiri
University of Bath
Bath, UK
vpn22@bath.ac.uk

C.V. Jawahar
IIIT Hyderabad
Hyderabad, Telengana
jawahar@iiit.ac.in

## ABSTRACT

Indian language machine translation performance is hampered due to the lack of large scale multi-lingual sentence aligned corpora and robust benchmarks. Through this paper, we provide and analyse an automated framework to obtain such a corpus for Indian language neural machine translation (NMT) systems. Our pipeline consists of a baseline NMT system, a retrieval module, and an alignment module that is used to work with publicly available websites such as press releases by the government. The main contribution towards this effort is to obtain an incremental method that uses the above pipeline to iteratively improve the size of the corpus as well as improve each of the components of our system. Through our work, we also evaluate the design choices such as the choice of pivoting language and the effect of iterative incremental increase in corpus size. Our work in addition to providing an automated framework also results in generating a relatively larger corpus as compared to existing corpora that are available for Indian languages. This corpus helps us obtain substantially improved results on the publicly available WAT evaluation benchmark and other standard evaluation benchmarks.

## CCS CONCEPTS

• **Information systems** → Data cleaning; Association rules; Presentation of retrieval results; **Structure and multilingual text search**; • **Computing methodologies** → Machine translation; *Information extraction*; • **Applied computing** → *Language translation*.

## KEYWORDS

parallel corpus, machine translation, information retrieval

---

[*]Both authors contributed equally to this work.

## 1 INTRODUCTION

Advances in machine translation, language-modelling, and other natural language-processing has led to a steep increase performance on tasks for many high-resource languages [Edunov et al. 2018; Ott et al. 2018]. One major driving factor is many western languages which become test-beds for the methods are already high-resource, which works in favour of methods which are data hungry [Koehn and Knowles 2017]. The high resource European counterparts have supporting projects like Europarl [Koehn 2005], Paracrawl [Bañón et al. 2020]. These have enabled large scale sentence aligned corpora to be developed. Similar efforts have not been realized for languages in the Indian subcontinent [Joshi et al. 2020]. Evidently, attempts need to be undertaken to improve this situation. Our work directly addresses this lacuna in the Indian machine translation setting. Specifically, through this work, we aim to achieve the following objectives:

- Provide a large scale sentence aligned corpus in 11 Indian languages, viz. CVIT-PIB corpus that is the largest multilingual corpus available for Indian languages as can be seen from Table 1.
- Demonstrate that such a corpus can be obtained automatically with no human effort using iterative alignment method.
- Provide robust standardized evaluation methodology and strong baselines that can be adopted and improved upon to ensure systematic progress in machine translation in Indian languages.

We briefly examine the alternatives to our approach and argue the need for adopting the proposed approach.

*Working at Scale.* There have been impressive works for low-resource languages at scale [Aharoni et al. 2019; Lepikhin et al. 2020; Schwenk et al. 2019], for instance working with 1620 language pairs [Schwenk et al. 2019] . However, not all these advances are feasible with regard to the compute resources available to standard academic research groups. Specifically, large models that converge faster, transfer more, and improve performance even for low-resource

languages [Aharoni et al. 2019; Lepikhin et al. 2020] are not trainable on hardware available to many research groups. Hence, we argue that this approach is not viable for Indian machine translation research, at this moment.

*Presently available corpora and baselines.* Unfortunately, research in Indian language machine translation suffers from the lack of suitable publicly available models and baselines. Those that are available are rather limited in scope or evaluation. For instance, a widely used corpus that is available is the ILCI Corpus [Jha 2010]. The corpus has 50K sentences aligned across many languages in the country. However, this corpus is limited to specific domains. The evaluation strategies on this corpora in literature also lacks comparability with no standard test-split. Despite these limitations, the corpus has been used by several reported sources as training data in literature to develop and study Machine Translation [Anthes 2010; Goyal et al. 2020; Kunchukuttan et al. 2014]. However, this corpus is not useful for applications like KR et al. [2019], due to the limitations. The Workshop on Asian Translation (WAT) [Nakazawa et al. 2019, 2017, 2018] on the other hand provides a standardized platform for a few languages. Similarly the Workshop on Machine Translation (WMT) [Barrault et al. 2019] from time to time hosts tasks with directions involving Indian Languages. Unfortunately, though several iterations of these tasks have concluded, to the best of our knowledge, there are no trained models that are publicly available at the moment. We summarize and list the presently available corpora in Table 1. It is evident that large scale multi-lingual corpora are lacking presently. There are multiple attempts in the past for developing ML solutions for Indian languages, including [Bhattacharyya et al. 2016] and other attempts such as [Goyal et al. 2020; Kunchukuttan and Bhattacharyya 2016; Murthy et al. 2019; Singh and Bhattacharyya 2019]. Unfortunately, most of these methods do not evaluate on a standard benchmarks/dataset for reliable comparison of the performance. They also have inferior performance to our approach and do not provide publicly available models. This thereby inhibits research in the community.

*Proposed approach.* We believe, that most methods applicable to the high-resource languages should work just as well in Indian languages in presence of the same amount of data. A simple solution which maintains Occam's razor is to change the low-resource situation, as more content is created online in many Indian Languages which is not pursued as much as it should be. Steps have been taken towards improving the situation in the monolingual corpus space [Kakwani et al. 2020; Kunchukuttan et al. 2020].

In this work we demonstrate how, using recent advances in Multilingual Neural Machine Translation (MNMT) [Dabre et al. 2020; Johnson et al. 2017; Vaswani et al. 2017] in an Expectation Maximization (EM) setup in the face of incomplete data, it is possible to change the status-quo of low resource to produce larger corpus and strong baselines in machine-translation for several Indian languages. A first-step towards this was taken in Siripragada et al. [2020] where we presented the CVIT-PIBv0.0 and CVIT Mann Ki Baat corpora. We substantially extend and refine this work through an iterative pipeline illustrated in Figure 1. We also co-opt some ideas proposed in low-resource adaptation for NMT proposed by
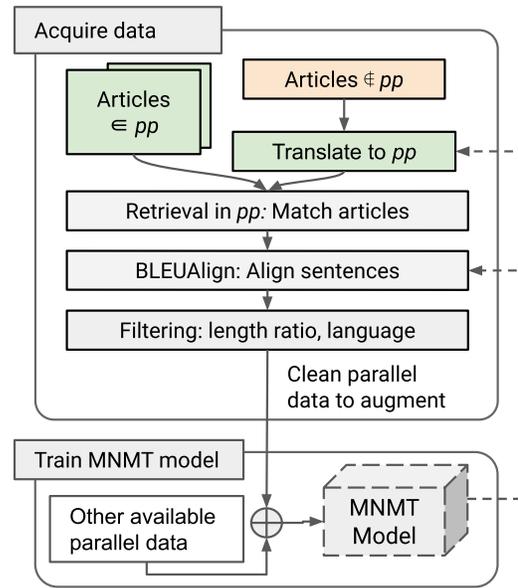


**Figure 1: Iterative alignment pipeline used for expanding the corpus for Indian languages. We observe that (i) A better MNMT model leads to better alignment and larger corpus (ii) Larger corpus leads to better MNMT model. We iterate until no further improvement is observed. The dashed lines indicate application of the trained MNMT model. *pp* stands for an arbitrary pivot language.**

[Sennrich and Zhang 2019]. In the process, we attain stronger baselines in translating the involved language-directions. Our contributions summarized are as follows:

(1) **Low Resource→High Resource**: We extensively study the iterative-alignment methods provided by Sennrich and Volk [2011] in the context of CVIT-PIBv0.0 dataset created for Siripragada et al. [2020]. Through successful execution of these methods, we increase the corpora size aggregated over all language-pairs from our previous 613K[1] to 2.78M (~ 353% increase) that we term the CVIT-PIB corpus v1.3.

(2) **Comparable and Strong Baselines**: We report consequent stronger baselines for MNMT in Indian languages from the improvement in data, validating the utility of the corpora we provide. The final MNMT model covers 11 languages and 110 language-directions with competitive or state-of-the-art performance in 12 tasks on public leaderboards.

(3) **Trained models and code**: We release the source-code, trained models and the datasets[2] to further research in this area and to aide applications that could be enabled by a functional MT. To the best of our knowledge, these are the only trained models available for translation with focus on Indian languages at the time of writing this document.

The rest of this document is organized as follows: In Section 2, we describe using traditional methods in MT based alignment to

---

[1]In Siripragada et al. [2020], we report this as 408K considering only English alignments, in this work we consider all-directions.
[2]http://preon.iiit.ac.in/~jerin/bhasha/

improve the parallel corpus across 11 Indian Languages. In Section 3 we report stronger baselines for MNMT in Indian languages across many available public tasks.

## 2 ITERATIVE ALIGNMENT FOR PIB

We apply the methods in Sennrich and Volk [2011] to iteratively improve parallel data. The procedure is analogous to expectation-maximization (EM) algorithm. In the expectation step, we use noisy alignments of parallel sentences from news articles to get a meaningful signal to obtain a better Maximum Likelihood Estimate (MLE) function for the MT model. In the maximization-step the improved MT model is used to obtain stronger alignments. Unlike Sennrich and Volk [2011], we use an MNMT model in place of the Statistical Machine Translation (SMT) model. In this section, we provide details about the corpus, the methods used to obtain the same and analysis of its characteristics.

### 2.1 Data Sources

To obtain an initial Multilingual NMT (MNMT) system, we rely on the datasets compiled from several sources listed in Table 1. We use backtranslation [Sennrich et al. 2016] to improve data in Hindi and Telugu.

The Press Information Bureau (PIB) is used in this work as a source for articles published in several Indian Languages to extract a multiparallel corpus. The PIB is very similar to a newspaper publishing in several languages except with strong one-to-one matches between documents and monotonic sentences which provide more parallel sentences through automatic sentence alignment algorithms. This section describes using the same crawled content as Siripragada et al. [2020] and focuses on improving the quality of alignments, in an attempt to consequently improve corpus size and performance of the MNMT model.

| Source | #pairs | #lang | type |
|---|---|---|---|
| IITB-en-hi [Kunchukuttan et al. 2017] | 1.5M | 2 | en-hi |
| UFAL EnTam [Ramasamy et al. 2012] | 170K | 2 | en-ta |
| WAT-ILMPC [Nakazawa et al. 2018] | 800K | 7 | xx-en |
| ILCI [Jha 2010] | 550K | 11 | xx-yy |
| OdiEnCorp [Parida et al. 2020] | 27K | 2 | en-or |
| Backtranslated-Hindi | 2.5M | 2 | en-hi |
| Backtranslated-Telugu | 500K | 2 | en-te |
| CVIT Mann Ki Baat[Siripragada et al. 2020] | 41K | 10 | xx-yy |
| PMIndia-Corpus[Haddow and Kirefu 2020] | 728K | 13 | xx-yy |
| CVIT-PIBv0.0[Siripragada et al. 2020] | 613K | 11 | xx-yy |
| CVIT-PIBv0.2 | 1.17M | 11 | xx-yy |
| **CVIT-PIBv1.3** | **2.78M** | 11 | xx-yy |

**Table 1: Publicly available corpuses for Indian languages. The last group of rows were not used for training. CVIT Mann Ki Baat is used for evaluation purposes only and has overlap with PMIndia Corpus. All other sources are used for training the multilingual model. xx-yy indicates parallel sentences aligned across multiple languages. Last row is the proposed corpus.**

We are aware of the existence of PMIndia [Haddow and Kirefu 2020], a source of similar nature and motivation as PIB, but make a conscious choice not to use it in this work to prevent data-leakage issues of possible overlap with one of our test-sets CVIT Mann Ki Baat [Siripragada et al. 2020].

### 2.2 Iterative Alignment

Our iterative alignment procedure requires document and sentence alignment algorithms, an MNMT formulation which is trained again with refreshed data in each iteration. We describe the constituent components illustrated in Figure 1 and describe the iterative procedure ahead.

#### 2.2.1 Text Processing.

*Text Cleaning and Standardization.* We allow for noise on the web in the pipeline and avoid any linguistic features. There is noise present in documents generated in the past with Indian languages content. A desideratum for retrieval and matching is that the model is capable of handling such noise. This could be unicode issues, non-standard or normalized text which is present all-across sources on the web. Some amount of noise is mitigated by past works [Bhattacharyya et al. 2016] by standardising unicode, scripts etc[3].

*Tokenization.* We use SentencePiece [Kudo and Richardson 2018] to tokenize sentences into subword-units. The subwords which cover a corpus are decided optimizing likelihood of a unigram language-model over a large corpus and candidate subwords in EM steps. Recent works [Edunov et al. 2018; Ng et al. 2019] addressing high-resource western-languages follow a joint vocabulary of some 32K-64K subwords as a subword model creation hyperparameter. We observed in our early experiments that in the presence of huge imbalance of data among languages and the huge difference in scripts unlike major European languages, for example, this approach leads to subwords which reduce to characters for the less represented languages. In order to avoid the artifacts from such a subword learning strategy, we instead choose 4K subwords for each of the languages involved and take a union of these to generate the final vocabulary[4]. The process results in a vocabulary of 40K subword-units tokens for 11 languages which we maintain fixed across all iterations. We note that the artifacts can also be mitigated by a temperature based sampling for sentences among languages as Aharoni et al. [2019].

*Filtering parallel pairs.* To obtain a filtered corpora at every iteration, we allow only sentence-pairs into to the training pipeline where source length to target length ratio is in [0.5, 2.0]. We also use langid[5], a language identifier through writing script to filter sentences with foreign language tokens.

#### 2.2.2 Alignment Algorithms.
To perform document alignment we translate the articles to a common pivot language. We use the translations to rank candidate-matches in the pivot language. SentencePiece tokenization of each sentence eliminates requirement of a curated stop-words list, enabling us to compute similarity in

---

[3]https://github.com/anoopkunchukuttan/indic_nlp_library
[4]This design choice is partially inspired by the reasoning and supporting experiments in Sennrich and Zhang [2019].
[5]https://github.com/saffsd/langid.py

the search space of any desired pivot language. Cosine similarity on the *term frequency - inverse document frequency* (tf-idf) [Buck and Koehn 2016] is put to use to rank retrieved articles in the space of pivot language. Search space to find and rank candidate articles matching a translation is restricted to only in a vicinity of dates (2 days) of posted news articles.

Upon obtaining aligned document pairs, we use Bleualign [Sennrich and Volk 2010], an MT based sentence alignment algorithm. Other conventional sentence length based alignment algorithms such as Gale-Church [Gale and Church 1993] also exist, but we rely on MT based alignment as the performance of the NMT model increases with every iteration resulting in better sentence alignment. Bleualign also aggressively filters reducing false matches [Bañón et al. 2020].

*2.2.3 Multilingual Neural Machine Translation (MNMT) Model.* We use fairseq [Ott et al. 2019] for training a Transformer-Base [Vaswani et al. 2017] based MNMT system. The model we begin with is same as our first multilingual model in Siripragada et al. [2020]. However, unlike Siripragada et al. [2020], to refine the CVIT-PIBv0.0 dataset further, we choose a many-to-English model formulation trained to translate only from non-English languages to English. This is advantageous because (1) it enables faster training and retraining time, (2) the setting provides more capacity to English decoding which improves translation performance and consequently – retrieval in English. The model parameters - encoder, decoder and embeddings are shared among all languages. We additionally use tied embeddings [Press and Wolf 2017] at the encoder and decoder. In this work, we denote our first many to many model with no CVIT-PIBv0.0 dataset augmentation as M2M-0, the following many to English models with incremental dataset variations as M2EN-1, M2EN-2, M2EN-3. We additionally consider the best model from Siripragada et al. [2020] after training with CVIT-PIBv0.0 dataset augmentation, denoted in this work as M2M-1, to attempt another translation to a non English and possibly related-language as an alternative for retrieval.

*2.2.4 Iterations.* In each iteration, we initialize training with the model from previous iteration (warm-start). This helps to reduce training time when compared to a model training from scratch (cold-start) as we benefit from learning in the previous iterations. To maintain a constant increment in articles, we set a threshold on keeping a retrieved candidate at a constant value for each language. We observe that the scores improve with successive iterations, consequently obtaining more matching documents. We stop the iteration process at a stage of diminishing returns, i.e there is no prospect of a justifiable increase in corpus size (See Figure 2). Note that in future, further expansion through an increasing number of documents is always possible.

## 2.3 Discussions

The many-languages involved and the disparity in sizes in training lead to a setting where we can dissect and study several aspects. First we study the iterative alignment process, comparing retrieval accuracy, BLEU scores and increase in corpora side by side.

We track BLEU [Papineni et al. 2002] for the MT model, a *pseudo* retrieval accuracy for the retrieval pipeline and count of successful
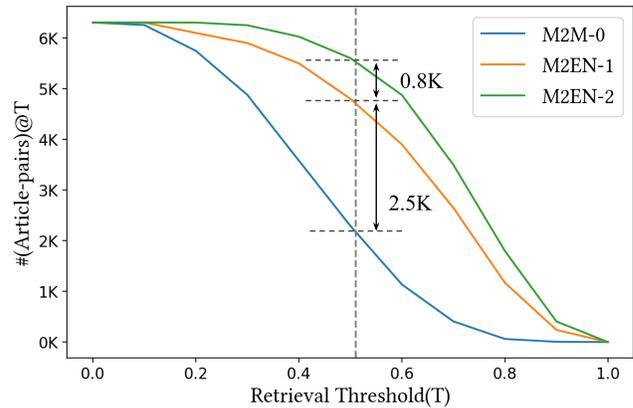


**Figure 2: Figure illustrating number of articles on the Y-axis obtained at a given threshold on X-axis. We maintain a constant threshold (dotted line) of *0.51* across model iterations M2M-0, M2EN-1 and M2EN-2. In case of Marathi, when compared to M2M-0 we acquire an additional 2.5K article pairs using M2EN-1 and 0.8K more using M2EN-2. From the graph, we observe saturation in articles pairs after iteration 2 indicating a point of diminishing returns.**

sentence alignments for increase in corpora over iterations. We employ the following technique to arrive at our *pseudo* retrieval accuracy. The articles which match in dates and ministry information from meta-information are collected. In many languages, there are enough true-positive matches to report the consequent evaluations as a proxy to retrieval accuracy.

These numbers reported in Table 2 indicate mostly consistent trends reflecting improvement in BLEU scores in translating to English for the CVIT Mann Ki Baat (WAT-2020 test split[6]), retrieval accuracies and consequently the resulting acquired data-sizes. The BLEU scores improve with the addition of more data, while the retrieval accuracy and data sizes improve with updating the parallel-corpus generated from PIB using the higher-performing MNMT model. Through successive iterations, the corpora increases in size and gets refined. We observe in Table 2, most increase for the languages with already good data (Hindi). In three iterations, we add a net 744K sentences aligned to English on top of the 408K (82% increase) sentences in the previous release, CVIT-PIBv0.0. But however, it is worth noting that some languages which are aided by the transfer and the new data have almost increased an order in sizes in iterations later (Marathi, Oriya and Bengali). Once the parallel corpus extraction from PIB is saturated (at M2EN-3), we crawl another 7 months of articles to expand the data further, and run another round of the alignment routine over the entire corpus. We obtain the multiparallel corpus by getting sentence alignments amongst other languages by bridging through the English part of the existing English-centric data. The process ends ups providing an additional 2.17M (353% increase) sentences to the previous 613K

---

[6]http://lotus.kuee.kyoto-u.ac.jp/WAT/indic-multilingual/

|  | Model | Languages | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  |  | hi | ta | ml | mr | gu | te | or | bn | pa |
| Retrieval Accuracy | M2M-0 | 76.77 | 54.12 | 45.63 | 34.05 | 52.52 | 24.06 | - | - | - |
|  | M2EN-1 | 86.91 | 71.36 | 69.77 | 47.57 | 63.67 | 62.3 | - | - | - |
|  | M2EN-2 | 92.39 | 80.84 | 80.4 | 51.89 | 75.86 | 70.05 | - | - | - |
| xx→en BLEU | M2M-0 | 17.63 | 9.49 | 11.2 | 11.35 | 14.44 | 6.98 | - | 10.82 | - |
|  | M2M-1 | 20.11 | 13.38 | 14.89 | 15.89 | 19.60 | 10.02 | - | 14.77 | - |
|  | M2EN-1 | 21.29 | 14.41 | 15.59 | 16.73 | 20.04 | 9.25 | - | 15.19 | - |
|  | M2EN-2 | 22.17 | 15.25 | 16.92 | 17.64 | 21.27 | 9.93 | - | 16.39 | - |
|  | M2EN-3 | 22.00 | 15.43 | 16.98 | 18.02 | 21.28 | 10.05 | - | 16.50 | - |
|  | M2EN-4 | 22.55 | 16.48 | 18.35 | 19.35 | 23.08 | 13.62 | - | 17.83 | - |
|  | M2EN-4-32K | 22.98 | 17.05 | 18.86 | 19.53 | 23.39 | 13.74 | - | 18.06 | - |
| Corpus Size | M2M-0 | 156.3K | 61.0K | 17.0K | 40.0K | 25.5K | 6.0K | 9.1K | 21.6K | 26.3K |
|  | M2EN-1 | 189.2K | 73.7K | 28.7K | 71.6K | 26.3K | 5.3K | 20.3K | 42.4K | 24.6K |
|  | M2EN-2 | 195.2K | 87.1K | 32K | 81.0K | 29.4K | 5.7K | 58.5K | 48.3K | 27.1K |
|  | Increment | 38.9K | 26.1K | 15.0K | 40.8K | 3.9K | 0 | 49.4K | 26.8K | 0.8K |

Table 2: Incremental improvements in Accuracy, *xx→en* BLEU scores on *Mann Ki Baat* (test-split from WAT-2020 used here to stay comparable) and Corpus size. We observe increments in retrieval accuracies consistent with increase in BLEU scores. WAT-2020 test split does not contain pa and or, while PIB does.

sentences aligned across languages, resulting in a corpora size of *2.78M*, viz. CVIT-PIBv1.3[7].

In the case of Oriya, we observe no date-based matches. This leads to depending entirely on accurate document pair retrieval for extracting a corpus of reasonable alignment accuracy. The preliminary efforts used to retrieve and align data in Siripragada et al. [2020] comprises of Bible [Parida et al. 2020] which has been observed to transfer poorly to other domains consequently to poor sentence level alignment accuracies for Oriya in our earlier models [Siripragada et al. 2020]. However, with every iteration in Tables 2 and 4 we observe increase in BLEU scores when translating to English. This leads to better retrieval and sentence level accuracies.

In our early experiments, we had left Urdu out from the iterative alignment procedures due to unresolved bugs in the text processing in the pipeline. This lead to no changes in Urdu to English corpus sizes for the first three iterations. However, this paves the way for a case study to evaluate the effect of the improvements in other multilingual language-pairs in translating to and from Urdu as we include Urdu in the performance evaluations. We notice that as the remaining resources improve and better alignments are put in place, and with no further Urdu data enhancements, there are improvements in the BLEU scores of language pairs involving Urdu, observed in Table 4 through M2EN-1 to M2EN-3.

To summarize, we provide a new corpus and a method to expand the corpus from publicly available sources.

- We make a new large sentence aligned corpus (CVIT-PIBv1.3) available for researchers, as a result of the iterative alignment described above. The detailed analysis of the iterative alignment procedure is provided in Table 2.
- This new corpus is possibly the largest sentence aligned multi-lingual corpus for Indian languages as can be observed through Table 1 in number.

---

[7]There exists a CVIT-PIBv0.2 used for WAT-2020 and WMT-2020. See Appendix for more information.

- Table 3 shows the corresponding increase in size and percentage increase from our previous effort using the same document set as in [Siripragada et al. 2020] with the increment obtained only through the proposed procedure. We further provide strong baselines that are presented in the next section.
- We also hint at a method that can allow continuous expansion of this corpus and eventually enabling a class of recent language processing methods on Indian languages.

Having described the process of iterative refining and enriching parallel-corpus resources for Indian Languages, we use the resulting corpora in training two models for obtaining baselines - one many-to-many (M2M) and the other many-to-English (M2EN). Since these are the next iterations, we label these M2M-4 and M2EN-4. In addition to these, we consider a model with 32K output vocabulary giving more vocabulary to English (M2M, M2EN uses 4K) which is denoted by M2EN-4-32K, disabling tied-embeddings. The three models are used to establish strong baselines for the 11 languages and consequent 110 directions involved which we describe in the next section.

## 3 STRONGER AND REPEATABLE BASELINES

It is important to further research to first take stock of where we are, and often simple baselines which compete in comparison to sophisticated methods serve this exact purpose [Arora et al. 2019; Xiao et al. 2018]. We consider the possibility using standard approaches which work for high-resource languages to provide such baselines for translation among Indian languages. Our previous work Philip et al. [2019] reports baselines for multilingual translation models for 5 Indian Languages and English. Further improvements are brought about in Siripragada et al. [2020] in several tasks adding more languages. It is also important to take care that these are "simple baselines" which are comparable [Post et al. 2012], leaving plenty of room for improvements ahead.

|    | en | hi | ta | te | ml | ur | bn | gu | mr | or | pa |
|----|----|----|----|----|----|----|----|----|----|----|----|
| en |    | 269594 | 118759 | 44888 | 44986 | 202578 | 93560 | 59739 | 117199 | 98230 | 103296 |
| hi |    |    | 64936 | 28562 | 27154 | 109946 | 49584 | 41583 | 69167 | 61065 | 75188 |
| ta |    |    |    | 17356 | 23599 | 48872 | 32988 | 29182 | 48527 | 44019 | 46340 |
| te |    |    |    |    | 10467 | 21141 | 17604 | 16325 | 18169 | 10462 | 25680 |
| ml |    |    |    |    |    | 20894 | 18136 | 18234 | 22793 | 19390 | 21960 |
| ur |    |    |    |    |    |    | 39290 | 29914 | 49683 | 43733 | 51817 |
| bn |    |    |    |    |    |    |    | 25154 | 34025 | 26449 | 35107 |
| gu |    |    |    |    |    |    |    |    | 30759 | 27140 | 35555 |
| mr |    |    |    |    |    |    |    |    |    | 46999 | 50411 |
| or |    |    |    |    |    |    |    |    |    |    | 43138 |

**Table 3: Multilingual shared content across language pairs for CVIT-PIBv1.3. Rows and columns indicate language pairs. The highlights are proportional to the increases in corpora sizes compared to the previous release CVIT-PIBv0.0 [Siripragada et al. 2020]**

Addressing simplicity, through the aforementioned works and this one, the model architecture and hence the learner's "capacity" is kept constant (Transformer-Base). Without increasing capacity, the experiments continue taking on all available translation tasks in Indian languages. A task here corresponds to language-directions or domains of datasets. The model we use is expected to be as general purpose as possible, after improving the data situation. There are no linguistics based priors in our methods or explicit handling of noise. This provides room for linguistics based improvements to build on simultaneously raising a call to revisit some older propositions. Addressing comparability, we comprehensively cover and compare with test-sets in prior-art and the shared tasks.

## 3.1 On Repeatability of Objective Evaluations

We address two important aspects (1) standard test-sets available and (2) reproducible evaluations.

*3.1.1 Test Sets.* We identify two class of test-sets among Indian languages, (1) which corresponds to ILCI which many early works evaluated translation quality on and (2) associated with the WAT or WMT tasks, which provide a leaderboard and standardized evaluations for comparison. To cover limitations of these test-sets, we propose a new test-set CVIT Mann Ki Baat in Siripragada et al. [2020]. We proceed to summarize how we compare to past work reporting numbers on these test-sets.

*3.1.2 Comparable reports of BLEU Scores.* Post [2018] addresses several issues of reproducibility and fair comparisons in reporting BLEU scores. In this work, we make our evaluations consistent with WAT leaderboard and provide a package to reproduce the procedure locally[8]. For WMT tasks, we report the values from the portal[9](gu-en) and the SacreBLEU[10] signature (ta-en). We make the hypotheses generated among all test-sets available[11] in case of a requirement to re-evaluate and compare with non-BLEU metric.

---

[8]https://github.com/jerinphilip/wateval
[9]http://matrix.statmt.org/matrix
[10]BLEU+case.mixed+numrefs.1+smooth.exp+tok.13a+version.1.4.12
[11]https://github.com/shashanksiripragada/generation-results

## 3.2 Results and Discussions

We begin by discussing the general merits of the model especially coming out of the multilingual formulations, and proceeding to elaborate on where we stand with regard to existing literature. For translating in to English, our best M2EN model (M2EN-4-32K) provide a cumulative improvement of +25 BLEU and an average improvement of +3.5 BLEU, compared to the previous best known multilingual model M2M-1 [Siripragada et al. 2020] with a similar coverage of languages, with the same model capacities on CVIT Mann Ki Baat test set in Table 2. This clearly points to the improvement consequent of change in data situation in the involved Indian languages. In Table 6, we report BLEU scores of the M2M-4 in a grid indicating the performance in the language-directions the model applies. We also highlight the improvement magnitude in color.

The M2EN model being stronger in the tasks it is trained towards compared to the corresponding M2M model, and is consistent with what is established in literature[Aharoni et al. 2019]. A reasoning for this disparity is not enabling temperature based sampling as in Aharoni et al. [2019] to balance out all language-pairs and correct for the imbalance in huge number of English aligned sentences existing in the training data. In Bapna and Firat [2019], corrections for the imbalance are observed to have led to degradation in high-resource languages. Note that both models are trained with the same capacity, and gain in BLEU scores and capability in translating more languages is a major advantage of the M2M model.

To study generalization, we take the models from the iterative procedure described in Section 2 and evaluate their performance on all available test-sets. The results can be observed in Tables 2 and 4.

Over the span of a few incremental works, we have significantly improved Hindi to English translation by a margin of +3 BLEU since Philip et al. [2019], obtaining higher numbers by using simple methods. This is unlike many rounds of distillation, hyperparameter optimization done by other groups to reach similar range of values[Nakazawa et al. 2018]. The other languages also have similar improvements in all directions.

Despite being not trained with the provided Gujarati data and using the data from ILCI and CVIT-PIBv1.3 corpus, we are able to achieve a BLEU score of 25.3 and 25.6 with our models M2EN-3 and M2M-3 respectively, competitive to the best BLEU score of 26.9 [Li

| direction | Model | IITB | UFAL | OdiEnCorp | WAT-ILMPC | | | | | | WMT19 | WMT20 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | hi | ta | or | hi | ta | te | ml | ur | bn | gu | ta (test) | ta (dev) |
| en→xx | M2M-0 | 19.83 | 6.78 | 4.29 | 19.99 | 10.86 | 16.19 | 7.19 | 13.27 | 9.69 | 6.4 | 3.5 | 4.8 |
| | M2M-1 | 20.52 | 7.31 | 5.26 | 20.39 | 11.63 | 16.63 | 7.92 | 16.55 | 9.58 | 9.5 | 4.2 | 6.0 |
| | M2M-3 | 21.20 | 7.22 | 4.78 | 20.92 | 11.95 | 17.10 | 7.70 | 15.78 | 10.13 | 11.3 | 4.9 | 7.1 |
| | M2M-4 | 21.28 | 7.80 | 5.25 | 20.25 | 10.00 | 15.80 | 6.60 | 16.08 | 9.29 | 12.5 | 5.1 | 7.4 |
| xx→en | M2M-0 | 21.94 | 18.64 | 11.05 | 27.99 | 17.82 | 21.63 | 11.97 | 20.57 | 16.67 | 17.9 | 12.9 | 12.7 |
| | M2M-1 | 22.48 | 19.76 | 10.84 | 28.31 | 18.65 | 22.58 | 12.71 | 21.16 | 16.77 | 23.6 | 14.3 | 13.9 |
| | M2M-3 | 23.07 | 19.87 | 12.07 | 28.99 | 19.16 | 23.96 | 12.77 | 21.15 | 17.38 | 25.6 | 15.9 | 15.3 |
| | M2M-4 | 22.84 | 19.66 | 12.28 | 27.88 | 18.09 | 22.93 | 12.19 | 21.19 | 16.74 | 25.2 | 16.6 | 15.1 |
| | M2EN-1 | 23.83 | 23.38 | 13.07 | 31.33 | 21.17 | 25.69 | 14.24 | 23.38 | 18.78 | 22.8 | 15.5 | 15.0 |
| | M2EN-2 | 24.65 | 25.32 | 15.62 | 32.88 | 23.19 | 28.11 | 15.68 | 24.53 | 20.03 | 24.5 | 16.6 | 16.3 |
| | M2EN-3 | 25.26 | 26.08 | 17.76 | 34.09 | 23.85 | 29.47 | 16.38 | 25.88 | 20.62 | 25.3 | 16.7 | 16.4 |
| | M2EN-4 | 25.01 | 26.49 | 17.41 | 33.73 | 23.35 | 30.14 | 15.87 | 26.38 | 19.89 | 24.6 | 17.2 | 16.6 |
| | M2EN-4-32K | 24.63 | 27.40 | 19.68 | 34.44 | 24.77 | 31.44 | 17.17 | 27.79 | 21.36 | 24.2 | 17.2 | 17.1 |

**Table 4: We report BLEU scores on available publicly available benchmark tasks for Indian Languages. The results on these benchmarks often have models that are specially tuned for various language pairs. We do observe that we obtain state of the art results on 4 of the language pairs and are competitive to other works that are more specific in most cases. This is despite not being specially tuned for these settings.**

et al. 2019] in WMT 2019 gu-en task. The best performing model did several rounds of backtranslation, distillation and multilingual formulation leveraging Hindi. Among these, we have only taken advantage of multilingual formulation and on top, pure data augmentation at the moment. A caveat here is that we have not put in efforts into filtering the test-data from the training data. But our corpus collection is independent of the news-sources WMT19 used and the Gujarati-English directions are as good as the claim, also supported by the results on CVIT Mann Ki Baat and ILCI test sets.

The BLEU scores on CVIT Mann Ki Baat test-set are provided in Table 6. We notice the most improvements in Oriya involved directions. Our overall multilingual model seems to have improved at M2M-4 in comparison to M2M-1. Despite not enforcing any linguistic priors, we get strong performance in many related languages.

It is also worth noting the correlations between Tables 3 and 6, that the highest improvements in BLEU has been for directions involving those languages for which more data has been added. With a model of fixed capacity throughout, simply increasing data has given increase in performance. The trend suggests the possibility to collect more sentence-aligned parallel text as a means to improve performance of machine translation models for Indian languages.

### 3.3 Comparison with Previous Works

*Non-standard comparisons.* Comparison is not standardized since previous methods evaluate on their own test set or on non-standard splits of the ILCI corpus that are not publicly available. This could be common in the initial stages of research for any community. Our work addresses it by evaluating it on more standard benchmarks with clear publicly available test splits. Kunchukuttan et al. [2014] attempts to build a collection of SMT models covering 11 Indian languages, similar to this work, except training and testing on splits from ILCI corpus (2K test-sentences, 500 for validation and remaining for training). However, the split is not available. Goyal et al. [2020] once again report numbers on ILCI, using a similar

| Work | IITB-hi-en | | WAT-ILMPC | |
|---|---|---|---|---|
| | en→hi | hi→en | en→hi | hi→en |
| SMT [Kunchukuttan et al. 2017] | 11.75 | 14.49 | - | - |
| NMT [Kunchukuttan et al. 2017] | 12.23 | 12.83 | - | - |
| Saini and Sahula [2018] | 18.22 | - | - | - |
| Philip et al. [2018] | **21.57** | 20.63 | - | - |
| Dabre et al. [2018] | | | 29.65 | 31.51 |
| Goyal and Sharma [2019] | - | 18.64 | - | - |
| Philip et al. [2019] | 20.17 | 22.62 | 26.25 | 31.55 |
| Siripragada et al. [2020] | 20.52 | 22.48 | 20.3 | 28.3 |
| Proposed Methods | 21.28 | **25.26** | 20.92 | **34.44** |

**Table 5: Comparison with publicly available baselines for English to Hindi and vice versa.**

split strategy as Kunchukuttan et al. [2014]. Murthy et al. [2019] compares with a similar test-set of 2K ILCI sentences. Similarly, Goyal et al. [2020] uses the ILCI test-set with a similar strategy to apply a formulation taking advantage of "related-languages". In our experiments we have found ILCI to be domain-specific (health, tourism) and providing false perception of high-scores by models which fail to generalize [Siripragada et al. 2020]. Due to the lack of reproducibility and comparability and the known demerits, we do not recommend future comparisons on any arbitrary ILCI test-set for benchmarking general purpose translation systems.

*Comparison on Hindi and English.* Hindi is the Indian language that has received highest attention with multiple attempts for translatiion to and from English. The results for comparison for Hindi-English on publicly available standard benchmarks that can be accessed are provided in Table 5. We obtain the highest scores for two of the tasks, i.e. IITB and WAT-ILMPC Hindi-English evaluations. Note that the highest BLEU score for IITB English to Hindi was obtained by our previous approach [Philip et al. 2018].

|      | bn    | en    | gu    | hi    | ml    | mr    | or     | ta    | te    | ur    | Δ      |
|------|-------|-------|-------|-------|-------|-------|--------|-------|-------|-------|--------|
| bn   |       | 16.79 | 15.50 | 21.62 | 5.75  | 11.05 | 12.42  | 4.99  | 5.64  | 24.73 | 34.36  |
| en   | 8.74  |       | 12.92 | 16.93 | 5.51  | 9.84  | 9.07   | 4.86  | 5.75  | 22.16 | 23.32  |
| gu   | 13.48 | 21.93 |       | 44.16 | 7.29  | 17.22 | 16.12  | 6.06  | 7.12  | 45.82 | 45.70  |
| hi   | 13.84 | 21.56 | 35.79 |       | 7.75  | 18.07 | 16.40  | 6.49  | 7.69  | 51.70 | 47.22  |
| ml   | 9.46  | 17.01 | 13.70 | 20.02 |       | 10.78 | 11.52  | 5.93  | 6.33  | 23.88 | 36.30  |
| mr   | 11.34 | 18.37 | 19.53 | 25.89 | 6.56  |       | 12.95  | 5.58  | 6.21  | 30.83 | 38.80  |
| or   | 12.98 | 19.36 | 19.94 | 26.99 | 6.21  | 13.19 |        | 4.96  | 5.65  | 26.92 | 71.06  |
| ta   | 8.32  | 15.30 | 11.51 | 17.20 | 5.80  | 9.26  | 10.04  |       | 5.70  | 20.56 | 33.31  |
| te   | 8.26  | 12.92 | 11.97 | 17.47 | 6.28  | 9.53  | 9.95   | 5.53  |       | 21.99 | 36.71  |
| ur   | 12.92 | 23.52 | 28.76 | 48.54 | 7.66  | 16.07 | 11.60  | 5.22  | 6.85  |       | 38.92  |
| Δ    | 24.12 | 24.16 | 38.96 | 38.92 | 13.47 | 31.19 | 109.16 | 11.74 | 22.32 | 91.66 |        |

**Table 6: BLEU scores of M2M-4 model on multilingual test set Mann-Ki Baat. Rows correspond to source languages and columns target languages. The colors indicate improvement (blue) or degradation (red) in comparison to M2M-1 [Siripragada et al. 2020]. We observe cumulative increment of +405 BLEU across all language pairs and a median increment of +2.7. The cumulative changes in translating to or from a given language in comparison to M2M-1 are provided under Δ header. It can be observed that related languages end up with higher BLEU scores without having to add the prior in the model formulation - e.g (hi, gu), (ur, gu), (ur, hi). Closely behind, there is (mr, hi) ahead of other language pairs.**

*Comparison on Public Leader Boards.* We provide detailed results for our method on publicly available leaderboards in Table 4 and Table 6. These can be used for comparisons and evaluations by various methods. As mentioned previously, these are on WAT tasks, WMT Tasks and CVIT Mann Ki Baat evaluation set. In all these tasks our models perform well obtaining state of the art results for several tasks. For instance, we obtain state of the art BLEU score of 19.68 on OdiEnCorp that is much higher than the previous state of the art of 8.6 and 24.77 BLEU score on WAT-ILMPC for Tamil to English that is higher than the previous state of the art 24.31.

We thus establish strong baselines for machine translation for Indian languages. Our multilingual model outperforms the previous works and even many carefully handcrafted MT systems for specific language pairs.

## 4 CONCLUSIONS AND FUTURE DIRECTIONS

In this work, we have made contributions to change the low-resource status of several Indian Languages for Machine Translation. Specifically (i) We introduced a large corpus that can enable Deep Machine Translation and associated research for these languages. (ii) Domain of these sentences allow to cover wider topics and practically more useful. More importantly, we established the utility of an algorithm that can help to grow the size starting from this state. More data will lead to better models, that in turn better alignment and more data. The corpus is bound to increase in size as more articles get added to PIB and the tooling in place to collect more sentences.

This work also possibly acts as the first NMT model dealing in Indian Languages that is publicly available for research. We hope this will spur more research within the research groups, specially in India. To enable the same, our code, models and data splits are made publicly available. Our corpus is now getting used in the WMT and WAT (International and Asian premier machine translation forums), demonstrating the utility. In addition to the parallel-corpus, we also make access to the crawling tools public - which can be used in the future to create document-level NMT datasets in Indian languages.

A challenging question will be the applicability of this method for online resources that are not really created by explicit translation. We believe, that a solution to this problem may not be too far from here. The methods used for search and alignment in this paper can be extended to use newspapers and news specific to a time-window in a weakly supervised setting with minimal human effort to enhance the parallel-corpus further. Using embeddings trained to mine parallel sentences have shown promise for High Resource Languages, which we will incorporate into this pipeline in the future. The meta-information on the stored PIB articles opens up possibilities to study document translation and active learning problems, left for future work.

Our M2EN models have high BLEU scores which allows for an application of backtranslation [Edunov et al. 2018; Sennrich et al. 2016] to improve the numbers further in the opposite direction (*en→xx*). Kim et al. [2020] reports scenarios where unsupervised NMT methods fail for the low-resource Gujarati-English pair due to limitations, and the enhancement of resources here implores a revisit.

With high-performing NMT systems to English from Indian languages, it is possible to create datasets and corpus for use in downstream tasks and by using the models provided by this work to further the research in Indian Languages, with Kunchukuttan et al. [2020] using the *CVIT Mann Ki Baat* to evaluate cross-lingual sentence retrieval being an example.

# REFERENCES

Roee Aharoni, Melvin Johnson, and Orhan Firat. 2019. Massively Multilingual Neural Machine Translation. *arXiv preprint arXiv:1903.00089* (2019).

Gary Anthes. 2010. Automated translation of indian languages. *Commun. ACM* 53, 1 (2010), 24–26.

Sanjeev Arora, Yingyu Liang, and Tengyu Ma. 2019. A simple but tough-to-beat baseline for sentence embeddings. In *5th International Conference on Learning Representations, ICLR 2017*.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-Scale Acquisition of Parallel Corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 4555–4567. https://www.aclweb.org/anthology/2020.acl-main.417

Ankur Bapna and Orhan Firat. 2019. Simple, Scalable Adaptation for Neural Machine Translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 1538–1548.

Loïc Barrault, Ondřej Bojar, Marta R Costa-Jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, et al. 2019. Findings of the 2019 conference on machine translation (wmt19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*. 1–61.

Pushpak Bhattacharyya, Mitesh M. Khapra, and Anoop Kunchukuttan. 2016. Statistical Machine Translation between Related Languages. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Tutorial Abstracts*. Association for Computational Linguistics, San Diego, California, 17–20. https://doi.org/10.18653/v1/N16-4006

Christian Buck and Philipp Koehn. 2016. Quick and Reliable Document Alignment via TF/IDF-weighted Cosine Distance. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*. Association for Computational Linguistics, Berlin, Germany, 672–678. https://doi.org/10.18653/v1/W16-2365

Raj Dabre, Chenhui Chu, and Anoop Kunchukuttan. 2020. A Comprehensive Survey of Multilingual Neural Machine Translation. *arXiv preprint arXiv:2001.01115* (2020).

Raj Dabre, Anoop Kunchukuttan, Atsushi Fujita, and Eiichiro Sumita. 2018. NICT's Participation in WAT 2018: Approaches Using Multilingualism and Recurrently Stacked Layers. In *Proceedings of the 32nd Pacific Asia Conference on Language, Information and Computation: 5th Workshop on Asian Translation: 5th Workshop on Asian Translation*. Association for Computational Linguistics, Hong Kong. https://www.aclweb.org/anthology/Y18-3003

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding Back-Translation at Scale. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. 489–500.

William A. Gale and Kenneth W. Church. 1993. A Program for Aligning Sentences in Bilingual Corpora. *Computational Linguistics* 19, 1 (1993), 75–102. https://www.aclweb.org/anthology/J93-1004

Vikrant Goyal, Sourav Kumar, and Dipti Misra Sharma. 2020. Efficient Neural Machine Translation for Low-Resource Languages via Exploiting Related Languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*. Association for Computational Linguistics, Online, 162–168. https://www.aclweb.org/anthology/2020.acl-srw.22

Vikrant Goyal and Dipti Misra Sharma. 2019. LTRC-MT Simple & Effective Hindi-English Neural Machine Translation Systems at WAT 2019. In *Proceedings of the 6th Workshop on Asian Translation*. Association for Computational Linguistics, Hong Kong, China, 137–140. https://doi.org/10.18653/v1/D19-5216

Barry Haddow and Faheem Kirefu. 2020. PMIndia–A Collection of Parallel Corpora of Languages of India. *arXiv preprint arXiv:2001.09907* (2020).

Girish Nath Jha. 2010. The TDIL Program and the Indian Langauge Corpora Intitiative (ILCI). In *Proceedings of the Seventh conference on International Language Resources and Evaluation (LREC'10)*.

Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, et al. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of ACL* (2017).

Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world. *arXiv preprint arXiv:2004.09095* (2020).

Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, Gokul N.C., Avik Bhattacharyya, Mitesh M. Khapra, and Pratyush Kumar. 2020. IndicNLPSuite: Monolingual Corpora, Evaluation Benchmarks and Pre-trained Multilingual Language Models for Indian Languages. In *Findings of EMNLP*.

Yunsu Kim, Miguel Graça, and Hermann Ney. 2020. When and Why is Unsupervised Neural Machine Translation Useless? *arXiv preprint arXiv:2004.10581* (2020).

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. Citeseer.

Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. *arXiv preprint arXiv:1706.03872* (2017).

Prajwal KR, Rudrabha Mukhopadhyay, Jerin Philip, Abhishek Jha, Vinay Namboodiri, and CV Jawahar. 2019. Towards automatic face-to-face translation. In *Proceedings of the 27th ACM International Conference on Multimedia*. 1428–1436.

Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226* (2018).

Anoop Kunchukuttan and Pushpak Bhattacharyya. 2016. Learning variable length units for SMT between related languages via Byte Pair Encoding. *arXiv preprint arXiv:1610.06510* (2016).

Anoop Kunchukuttan and Pushpak Bhattacharyya. 2020. Utilizing Language Relatedness to improve Machine Translation: A Case Study on Languages of the Indian Subcontinent. *arXiv preprint arXiv:2003.08925* (2020).

Anoop Kunchukuttan, Divyanshu Kakwani, Satish Golla, Avik Bhattacharyya, Mitesh M Khapra, Pratyush Kumar, et al. 2020. AI4Bharat-IndicNLP Corpus: Monolingual Corpora and Word Embeddings for Indic Languages. *arXiv preprint arXiv:2005.00085* (2020).

Anoop Kunchukuttan, Pratik Mehta, and Pushpak Bhattacharyya. 2017. The iit bombay english-hindi parallel corpus. *arXiv preprint arXiv:1710.02855* (2017).

Anoop Kunchukuttan, Abhijit Mishra, Rajen Chatterjee, Ritesh Shah, and Pushpak Bhattacharyya. 2014. Shata-Anuvadak: Tackling Multiway Translation of Indian Languages. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*. 1781–1787.

Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2020. GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding. *arXiv preprint arXiv:2006.16668* (2020).

Bei Li, Yinqiao Li, Chen Xu, Ye Lin, Jiqiang Liu, Hui Liu, Ziyang Wang, Yuhao Zhang, Nuo Xu, Zeyang Wang, et al. 2019. The NiuTrans machine translation systems for WMT19. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*. 257–266.

Rudra Murthy, Anoop Kunchukuttan, and Pushpak Bhattacharyya. 2019. Addressing word-order Divergence in Multilingual Neural Machine Translation for extremely Low Resource Languages. Association for Computational Linguistics, Minneapolis, Minnesota, 3868–3873. https://doi.org/10.18653/v1/N19-1387

Toshiaki Nakazawa, Nobushige Doi, Shohei Higashiyama, Chenchen Ding, Raj Dabre, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, Shantipriya Parida, et al. 2019. Overview of the 6th workshop on Asian translation. In *Proceedings of the 6th Workshop on Asian Translation*. 1–35.

Toshiaki Nakazawa, Shohei Higashiyama, Chenchen Ding, Hideya Mino, Isao Goto, Hideto Kazawa, Yusuke Oda, Graham Neubig, and Sadao Kurohashi. 2017. Overview of the 4th Workshop on Asian Translation. In *Proceedings of the 4th Workshop on Asian Translation (WAT2017)*. 1–54.

Toshiaki Nakazawa, Katsuhito Sudoh, Shohei Higashiyama, Chenchen Ding, Raj Dabre, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, and Sadao Kurohashi. 2018. Overview of the 5th Workshop on Asian Translation. In *Proceedings of the 32nd Pacific Asia Conference on Language, Information and Computation: 5th Workshop on Asian Translation: 5th Workshop on Asian Translation*. Association for Computational Linguistics, Hong Kong. https://www.aclweb.org/anthology/Y18-3001

Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, and Sergey Edunov. 2019. Facebook FAIR's WMT19 News Translation Task Submission. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*. 314–319.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A Fast, Extensible Toolkit for Sequence Modeling. In *NAACL (Demonstrations)*.

Myle Ott, Sergey Edunov, David Grangier, and Michael Auli. 2018. Scaling Neural Machine Translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*. Association for Computational Linguistics, Brussels, Belgium.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics, 311–318.

Shantipriya Parida, Ondřej Bojar, and Satya Ranjan Dash. 2020. OdiEnCorp: Odia–English and Odia-Only Corpus for Machine Translation. In *Smart Intelligent Computing and Applications*. Springer, 495–504.

Jerin Philip, Vinay P. Namboodiri, and C.V. Jawahar. 2018. CVIT-MT Systems for WAT-2018. In *Proceedings of the 32nd Pacific Asia Conference on Language, Information and Computation: 5th Workshop on Asian Translation: 5th Workshop on Asian Translation*. Association for Computational Linguistics, Hong Kong. https://www.aclweb.org/anthology/Y18-3010

Jerin Philip, Vinay P Namboodiri, and CV Jawahar. 2019. A baseline neural machine translation system for indian languages. *arXiv preprint arXiv:1907.12437* (2019).

Martin Popel and Ondřej Bojar. 2018. Training tips for the transformer model. *The Prague Bulletin of Mathematical Linguistics* 110, 1 (2018), 43–70.

Matt Post. 2018. A call for clarity in reporting BLEU scores. *arXiv preprint arXiv:1804.08771* (2018).

Matt Post, Chris Callison-Burch, and Miles Osborne. 2012. Constructing parallel corpora for six indian languages via crowdsourcing. In *Proceedings of the Seventh Workshop on Statistical Machine Translation*. Association for Computational Linguistics, 401–409.

Ofir Press and Lior Wolf. 2017. Using the Output Embedding to Improve Language Models. *EACL 2017* (2017), 157.

Loganathan Ramasamy, Ondřej Bojar, and Zdeněk Žabokrtský. 2012. Morphological Processing for English-Tamil Statistical Machine Translation. In *Proceedings of the Workshop on Machine Translation and Parsing in Indian Languages (MTPIL-2012)*. 113–122.

Sandeep Saini and Vineet Sahula. 2018. Neural machine translation for English to Hindi. In *2018 Fourth International Conference on Information Retrieval and Knowledge Management (CAMP)*. IEEE, 1–6.

Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2019. Wikimatrix: Mining 135m parallel sentences in 1620 language pairs from wikipedia. *arXiv preprint arXiv:1907.05791* (2019).

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving Neural Machine Translation Models with Monolingual Data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 86–96.

Rico Sennrich and Martin Volk. 2010. MT-based sentence alignment for OCR-generated parallel texts. In *The Ninth Conference of the Association for Machine Translation in the Americas (AMTA 2010)*.

Rico Sennrich and Martin Volk. 2011. Iterative, MT-based sentence alignment of parallel texts. In *Proceedings of the 18th Nordic Conference of Computational Linguistics (NODALIDA 2011)*. 175–182.

Rico Sennrich and Biao Zhang. 2019. Revisiting Low-Resource Neural Machine Translation: A Case Study. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 211–221. https://doi.org/10.18653/v1/P19-1021

Karanveer Singh and Pushpak Bhattacharyya. 2019. NMT in Low Resource Scenario : A Case Study in Indian Languages.

Shashank Siripragada, Jerin Philip, Vinay P. Namboodiri, and C V Jawahar. 2020. A Multilingual Parallel Corpora Collection Effort for Indian Languages. In *Proceedings of The 12th Language Resources and Evaluation Conference*. European Language Resources Association, Marseille, France, 3743–3751. https://www.aclweb.org/anthology/2020.lrec-1.462

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.

Bin Xiao, Haiping Wu, and Yichen Wei. 2018. Simple baselines for human pose estimation and tracking. In *Proceedings of the European conference on computer vision (ECCV)*. 466–481.

## A  RELATED LANGUAGES IN RETRIEVAL

The languages Marathi (mr), Gujarati (gu), Punjabi (pa) are similar to Hindi and exhibit high BLEU scores (Table 6) when translated to Hindi. They are also known to be similar [Kunchukuttan and Bhattacharyya 2020], so we experiment with hi as a pivot language. However, we found poor retrieval performance when compared to pivoting through English. Upon closer inspection, we observed that Hindi articles are much more elaborate while describing content while the mr, gu, pa equivalents are often summarized. This is evident when considering examples of Gujarati articles, as PIB offices of Gujarat are responsible for posting the articles in Gujarati and their respective English translations. We illustrate this phenomenon through an example in Figure 3, where we observe higher retrieval scores overall when compared to Hindi-based retrieval. The above analysis points to the success and a potential use-case of our model in being able to deliver consistent content across all languages for websites like PIB.
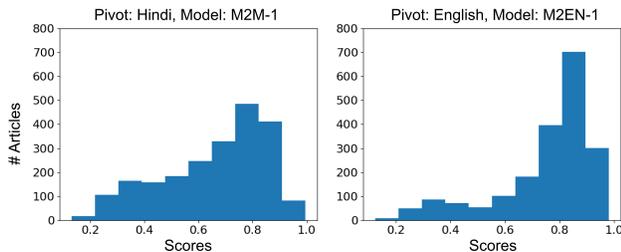


**Figure 3: Retrieval scores of *Gujarati*. Left bar chart indicates retrieval scores in case of model M2M-1 and pivot language *hi*. Right chart indicates scores in case of M2EN-1 and pivot *en*.**

## B  IMPLEMENTATION CONSIDERATIONS

The training was done on a machine equipped with 4 x NVIDIA 1080Tis or 2080Tis, 40 or 20 CPU cores and 128GB memory, depending on the allocation in our cluster. The multiparallel nature of the evolving PIB dataset and the ILCI dataset leads to an $O(N^2)$ growth with increase in samples across language pairs, specifically to train an M2M model. The training of the M2M models took ~3 days and an M2EN model took ~1 day from scratch, for the same number of epochs. Our hardware could not train the Transformer-Big [Vaswani et al. 2017] with the stability techniques prescribed by [Popel and Bojar 2018], so we had to use the Transformer-Base model.

We describe the code and modifications we implemented to train these models. We use fairseq [Ott et al. 2019] framework with some modifications[12]. Our modifications for this work include a dataloader equipped with memory-mapped storage using LMDB for fast access from the corpora described in Table 1. The non-standard SentencePiece routine required some additional integrations, and

[12]https://github.com/jerinphilip/fairseq-ilmt

these are publicly available as well[13]. We provide wrappers for using our trained models for inference in Python packaging, with models available for download separately.

| Device | Sentences | Batched | Time |
|--------|-----------|---------|------|
| CPU | 100 | Yes | 24.08 |
| GPU | 100 | Yes | 2.06 |
| CPU | 100 | No | 60.84 |
| GPU | 100 | No | 30.69 |

**Table 7: Caption**

We benchmarked the inference pipeline on both CPU and GPU machines. We present the summary of time taken to translate a sample test-set of sentences in Table 7. The inference can be done on a CPU for some practical use-cases, and our hosted demo model[14] runs on a CPU.

## C  ILCI NUMBERS

In Table 8, we provide comparisons with ILCI, with Kunchukuttan et al. [2014]. Despite not being domain adapted to ILCI, we obtain better BLEU scores in a majority of language-directions.

ILCI has been used by several works in the past with non-standard or non-reproducible performance benchmarks, which has rendered comparison hard with these. We urge the community to avoid using splits in ILCI for publishing results, absent any method to reproduce the constituent sentences in the split.

## D  VERSIONING

Table 9 provides the sizes of CVIT PIBv0.2, the corpus used in WMT-2020 and WAT-2020. The corpus is generated with M2EN-2 using articles crawled from PIB posted until December 2019.

CVIT PIBv1.3 contains articles crawled until August 2020. Future releases will be described on the project website[15].

[13]https://github.com/jerinphilip/ilmulti
[14]http://preon.iiit.ac.in/babel/gui
[15]http://preon.iiit.ac.in/~jerin/bhasha page.

|    | bn | en | gu | hi | ml | mr | pa | ta | te | ur |
|----|----|----|----|----|----|----|----|----|----|----|
| bn |    | 22.13 | 26.34 | 31.72 | 8.28 | 18.62 | 24.45 | 7.12 | 13.19 | 26.10 |
| en | 15.13 |    | 24.82 | 30.42 | 6.10 | 17.12 | 23.05 | 5.52 | 10.44 | 24.29 |
| gu | 22.91 | 30.32 |    | 54.47 | 8.88 | 28.38 | 42.12 | 8.38 | 15.64 | 45.56 |
| hi | 24.26 | 31.19 | 53.20 |    | 9.44 | 32.06 | 54.55 | 9.61 | 17.67 | 60.13 |
| ml | 14.34 | 18.01 | 20.13 | 24.49 |    | 14.04 | 18.88 | 6.41 | 10.82 | 20.88 |
| mr | 20.70 | 26.83 | 37.12 | 43.62 | 8.50 |    | 33.00 | 7.47 | 14.11 | 35.20 |
| pa | 22.42 | 31.47 | 49.28 | 68.46 | 9.21 | 29.17 |    | 8.92 | 16.27 | 54.94 |
| ta | 10.77 | 14.36 | 16.53 | 20.82 | 5.76 | 11.22 | 16.63 |    | 8.60 | 17.65 |
| te | 16.02 | 21.04 | 25.90 | 29.90 | 7.12 | 16.44 | 23.34 | 6.80 |    | 25.01 |
| ur | 19.80 | 28.18 | 42.60 | 56.41 | 8.54 | 24.33 | 43.23 | 8.70 | 14.39 |    |

**Table 8: BLEU scores of inference of model M2M-3 on random test split from ILCI. The reds indicate poorer performance compared to Kunchukuttan et al. [2014] and the blues better performance. Overall, our model performs better in 52/90 tasks and is +121 BLEU points ahead of *Sata Anuvaadak* with a median BLEU increase of 1.1.**

|    | en | hi | ta | te | ml | ur | bn | gu | mr | or | pa |
|----|----|----|----|----|----|----|----|----|----|----|----|
| en | 195208 | 87113 | 5752 | 31974 | 45344 | 48354 | 29421 | 80760 | 58461 | 27117 | |
| hi |    | 44031 | 3083 | 17819 | 11695 | 24849 | 19730 | 45950 | 36317 | 11442 | |
| ta |    |    | 3218 | 15029 | 4964 | 19175 | 16934 | 33636 | 27668 | 9150 | |
| te |    |    |    | 2543 | 415 | 1883 | 2625 | 2627 | 1834 | 1220 | |
| ml |    |    |    |    | 2378 | 9940 | 10132 | 14474 | 9843 | 4961 | |
| ur |    |    |    |    |    | 3795 | 2397 | 4941 | 3209 | 5584 | |
| bn |    |    |    |    |    |    | 10554 | 19914 | 14606 | 5332 | |
| gu |    |    |    |    |    |    |    | 17169 | 13682 | 5581 | |
| mr |    |    |    |    |    |    |    |    | 31377 | 9601 | |
| or |    |    |    |    |    |    |    |    |    | 6813 | |

**Table 9: Multilingual shared content across language pairs for CVIT-PIBv0.2. Rows and columns indicate language pairs. The highlights are proportional to the change after the iterative alignment process, reds indicating decrease and blues indicating increases in corpora sizes compared to the previous release v0.0. Evident from the table, we notice major increments in Marathi, Oriya and other languages. Tamil and Hindi which we had enough to be considered mid-to-high resource gain significant number as well. The maximum decrease is -283 for Telugu, which is negligible compared to the improvements of the order of ten-thousands in many language-pairs.**