

BitextEdit: Automatic Bitext Editing for Improved Low-Resource Machine Translation

Eleftheria Briakou, Sida I. Wang, Luke Zettlemoyer, Marjan Ghazvininejad
ebriakou@umd.edu, sida@fb.com, lsz@fb.com, ghazvini@fb.com



Bitexts are not Always Parallel

He was born in London.

त्याचा जन्म लंडनमध्ये झाला.

GLOSS: He was born in London.

Bitexts are not Always Parallel

She visited her sister.

ते डॉक्टरांना भेट देत आहेत.

GLOSS: They are visiting the doctor.

He was born in London.

त्याचा जन्म लंडनमध्ये झाला.

GLOSS: He was born in London.

Bitexts are not Always Parallel

She visited her sister.

ते डॉक्टरांना भेट देत आहेत.

GLOSS: They are visiting the doctor.

He was born in London.

त्याचा जन्म लंडनमध्ये झाला.

GLOSS: He was born in London.

I am not going back.

मांजर तिचे अन्न खात आहे.

GLOSS: The cat is eating her food.

Audits of Mined Bitext Reveal Systematic Quality Issues

Quality at a Glance: An Audit of Web-Crawled Multilingual Datasets

Julia Kreutzer^{a,b}, Isaac Caswell^a, Lisa Wang^a, Ahsan Wahab^c, Daan van Esch^a, Nasanbayar Ulzii-Orshikh^d, Allahsera Tapo^{b,c}, Nishant Subramani^{b,d}, Artem Sokolov^a, Claytone Sikasote^{b,d}, Monang Setyawan^b, Supheakmungkol Sarin^b, Sokhar Samb^{b,d}, Benoît Sagot^f, Clara Rivera^g, Annette Rios^h, Isabel Papadimitriouⁱ, Salomey Osei^{b,m}, Pedro Ortiz Suarez^{j,n}, Iroro Orife^{b,o}, Kelechi Ogueji^{b,p}, Andre Niyongabo Rubungu^{b,q}, Toan Q. Nguyen^r, Mathias Müller^k, André Müller^h, Shamsuddeen Hassan Muhammad^{b,s}, Nanda Muhammad^b, Ayanda Mnyakeni^b, Jamshidbek Mirzakhlov^{c,l}, Tapiwanashe Matangira^b, Colin Leong^b, Nze Lawson^b, Sneha Kudugunta^a, Yacine Jernite^{b,u}, Mathias Jenny^b, Orhan Firat^{b,c}, Bonaventure F. P. Dossou^{b,v}, Sakhile Dlamini^b, Nisansa de Silva^w, Sakine Çabuk Ballı^h, Stella Biderman^z, Alessia Battisti^h, Ahmed Baruwa^{b,y}, Ankur Bapna^a, Pallavi Baljekar^a, Israel Abebe Azime^{b,z}, Ayodele Awokoya^{b,z}, Duygu Ataman^{c,k}, Orevaoghene Ahia^{b,o}, Oghenefego Ahia^b, Sweta Agrawal^l, Mofetoluwa Adeyemi^{b,γ},

^aGoogle Research, ^bMasakhane NLP, ^cTurkic Interlingua, ^dHaverford College, ^eRobotsMali, ^fIntel Labs, ^gUniversity of Zambia, ^hGoogle, ⁱAIMS-AMMI, ^jInria, ^kUniversity of Zurich, ^lStanford University,

^mKwame Nkrumah University of Science and Technology, ⁿSorbonne Université, ^oNiger-Volta LTI, ^pUniversity of Waterloo, ^qUniversity of Electronic Science and Technology of China, ^rUniversity of Notre Dame, ^sBayero University Kano, ^tUniversity of South Florida, ^uHugging Face, ^vJacobs University Bremen, ^wUniversity of Moratuwa, ^xEleutherAI, ^yObafemi Awolowo University, ^zUniversity of Ibadan, ^{aa}Instadeep, ^{ab}University of Maryland, ^{ac}Defence Space Administration Abuja, ^{ad}Allen Institute for Artificial Intelligence

Abstract

With the success of large-scale pre-training and multilingual modeling in Natural Language Processing (NLP), recent years have seen a proliferation of large, web-mined text datasets covering hundreds of languages. We manually audit the quality of 205 language-specific corpora released with five major public datasets (CCAligned, ParaCrawl, WikiMatrix, OSCAR, mC4). Lower-resource corpora have systematic issues: At least 15 corpora have no usable text, and a significant fraction contains less than 50% sentences of acceptable quality. In addition, many are mislabeled or use non-standard/ambiguous language codes. We demonstrate that these issues are easy to detect even for non-proficient speakers, and supplement the human audit with automatic analyses. Finally, we recommend techniques to evaluate and improve multilingual corpora and discuss potential risks that come with low-quality data releases.

1 Introduction

Access to multilingual datasets for NLP research has vastly improved over the past years. A variety of web-derived collections for hundreds of languages is available for anyone to download, such as ParaCrawl (Esplà et al., 2019; Bañón et al., 2020), WikiMatrix (Schwenk et al., 2021) CCAAligned (El-Kishky et al., 2020), OSCAR (Ortiz Suárez et al., 2019; Ortiz Suárez et al., 2020), and several others. These have in turn enabled a variety of highly multilingual models, like mT5 (Xue et al., 2021), M2M-100 (Fan et al., 2020), M4 (Arivazhagan et al., 2019).

Curating such datasets relies on the websites giving clues about the language of their contents (e.g. a language identifier in the URL) and on automatic language classification (LangID). It is commonly known that these automatically crawled and filtered datasets tend to have overall lower quality than hand-curated collec-

The vast majority of low-resource languages contain less than 50% valid translations.

[Kreutzer et al.]

Detecting Fine-Grained Cross-Lingual Semantic Divergences without Supervision by Learning to Rank

Eleftheria Briakou and Marine Carpuat

Department of Computer Science
University of Maryland
College Park, MD 20742, USA
ebriakou@cs.umd.edu, marine@cs.umd.edu

Abstract

Detecting fine-grained differences in content conveyed in different languages matters for cross-lingual NLP and multilingual corpora analysis, but it is a challenging machine learning problem since annotation is expensive and hard to scale. This work improves the prediction and annotation of fine-grained semantic divergences. We introduce a training strategy for multilingual BERT models by learning to rank synthetic divergent examples of varying granularity. We evaluate our models on the Rationalized English-French Semantic Divergences, a new dataset released with this work, consisting of English-French sentence-pairs annotated with semantic divergence classes and token-level rationales. Learning to rank helps detect fine-grained sentence-level divergences more accurately than a strong sentence-level similarity model, while token-level predictions have the potential of further distinguishing between coarse and fine-grained divergences.

1 Introduction

Comparing and contrasting the meaning of text conveyed in different languages is a fundamental NLP task. It can be used to curate clean parallel corpora for downstream tasks such as machine translation (Koehn et al., 2018), cross-lingual transfer learning, or semantic modeling (Ganitkevitch et al., 2013; Conneau and Lample, 2019), and it is also useful to directly analyze multilingual corpora. For instance, detecting the commonalities and divergences between sentences drawn from English and French Wikipedia articles about the same topic would help analyze language bias (Bao et al., 2012; Massa and Scrinzi, 2012), or mitigate differences in coverage and usage across languages (Yeung et al., 2011; Wulczyn et al., 2016; Lemmerich et al., 2019). This requires not only detecting coarse content mismatches, but also fine-grained differences

in sentences that overlap in content. Consider the following English and French sentences, sampled from the WikiMatrix parallel corpus. While they share important content, highlighted words convey meaning missing from the other language:

EN *Alexander Muir's "The Maple Leaf Forever" served for many years as an unofficial Canadian national anthem.*

FR *Alexander Muir compose The Maple Leaf Forever (en) qui est un chant patriotique pro canadien anglais.*

GLOSS *Alexander Muir composes The Maple Leaf Forever which is an English Canadian patriotic song.*

We show that explicitly considering diverse types of semantic divergences in bilingual text benefits both the annotation and prediction of cross-lingual semantic divergences. We create and release the Rationalized English-French Semantic Divergences corpus (REFRED), based on a novel divergence annotation protocol that exploits rationales to improve annotator agreement. We introduce Divergent mBERT, a BERT-based model that detects fine-grained semantic divergences without supervision by learning to rank synthetic divergences of varying granularity. Experiments on REFRED show that our model distinguishes semantically equivalent from divergent examples much better than a strong sentence similarity baseline and that unsupervised token-level divergence tagging offers promise to refine distinctions among divergent instances. We make our code and data publicly available.¹

¹Implementations of Divergent mBERT can be found at: <https://github.com/Elbriakling-SemDiv>; the REFRED dataset is hosted at: <https://github.com/Elbriakling-SemDiv/tree/master/REFRED>.

Only 36% of En-Fr WikiMatrix sample are exact translations
[Briakou & Carpuat]

BitextEdit: Automatic Bitext Editing for Improved Low-Resource Machine Translation.

Eleftheria Briakou, Sida I. Wang, Luke Zettlemoyer and Marjan Ghazvininejad. NAACL Findings 2022.

Bitext Filtering as the Standard Approach to Bitext Quality Improvement

She visited her sister.

ते डॉक्टरांना भेट देत आहेत.

He was born in London.

त्याचा जन्म लंडनमध्ये झाला.

I am not going back.

मांजर तिचे अन्न खात आहे.

Bitext Filtering as the Standard Approach to Bitext Quality Improvement

She visited her sister.

ते डॉक्टरांना भेट देत आहेत.

He was born in London.

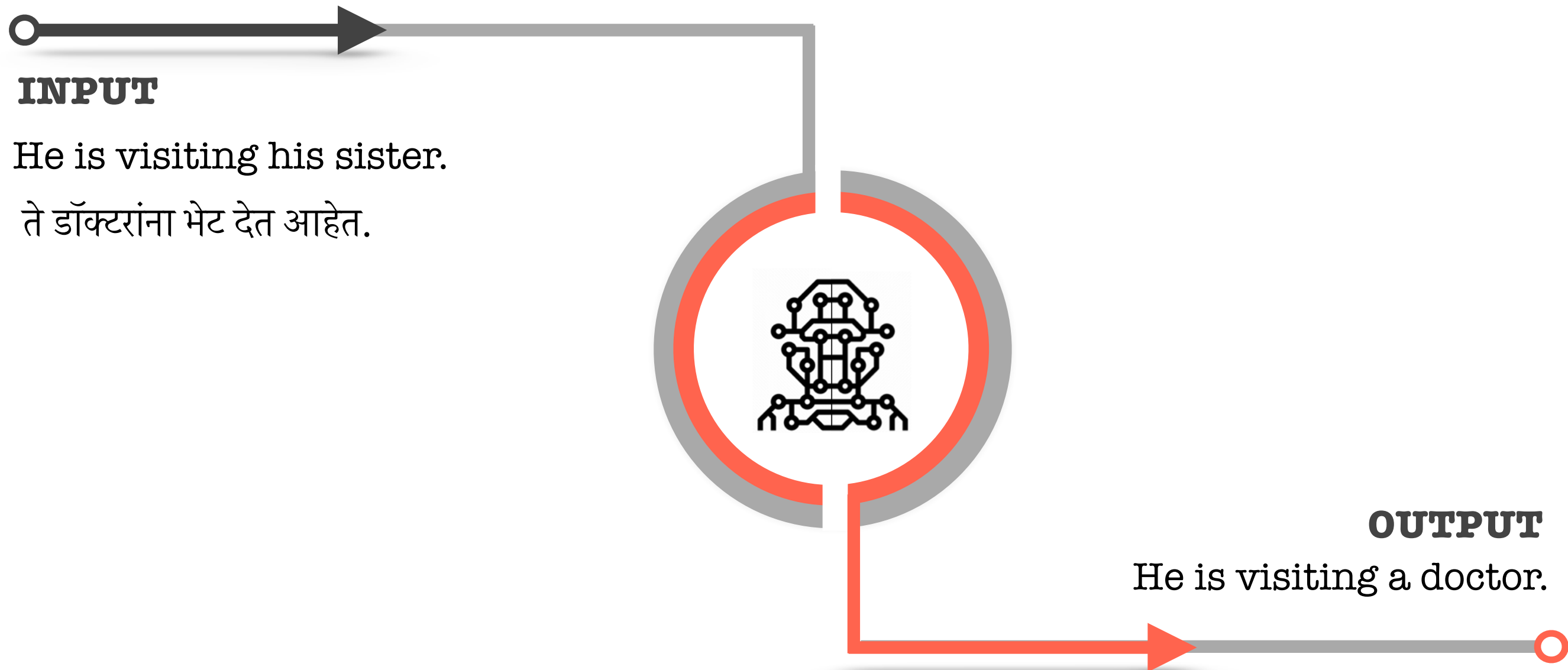
त्याचा जन्म लंडनमध्ये झाला.

I am not going back.

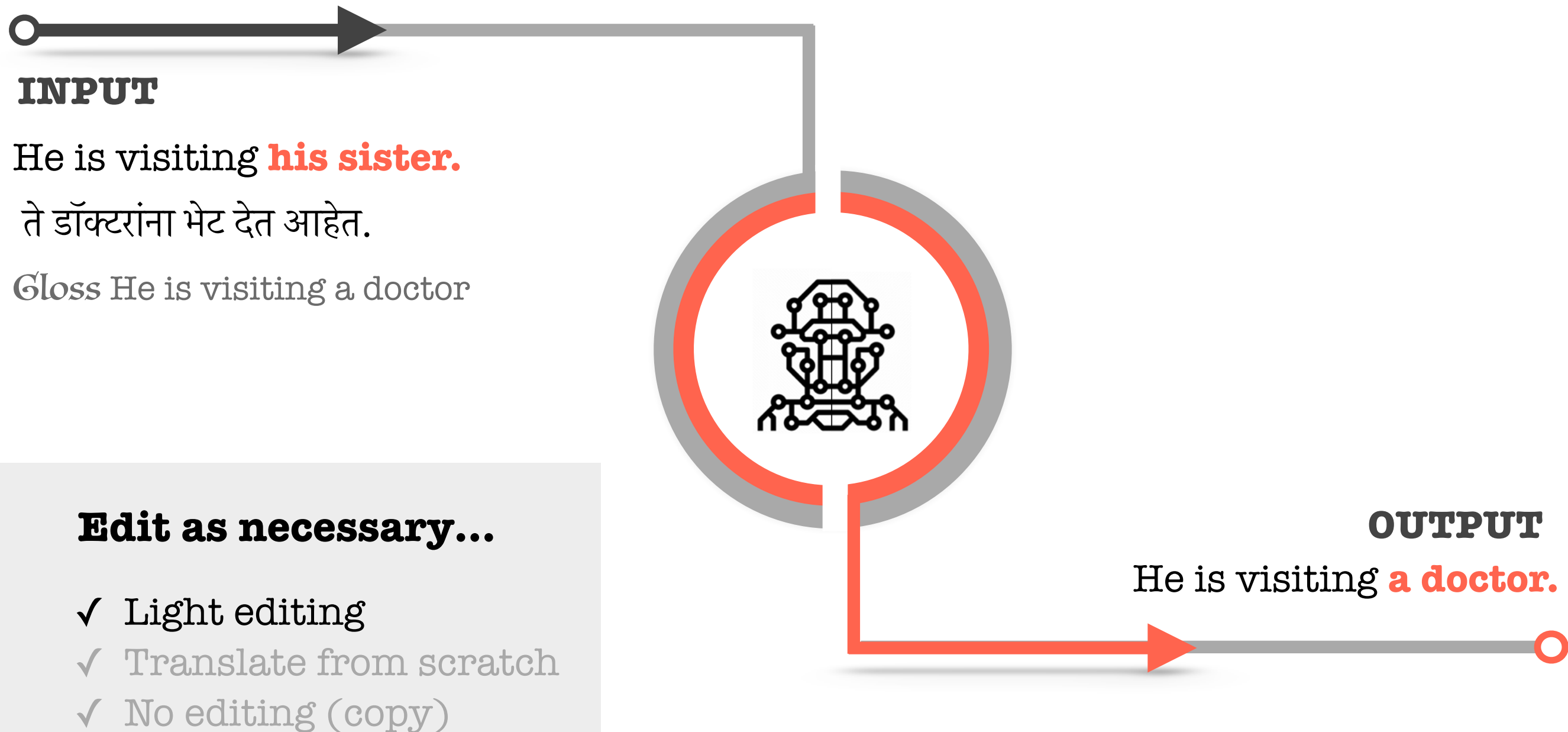
मांजर तिचे अन्न खात आहे.

- What if we cannot afford filtering (e.g., low-resource)?
- How do we handle imperfect translations beyond noise?

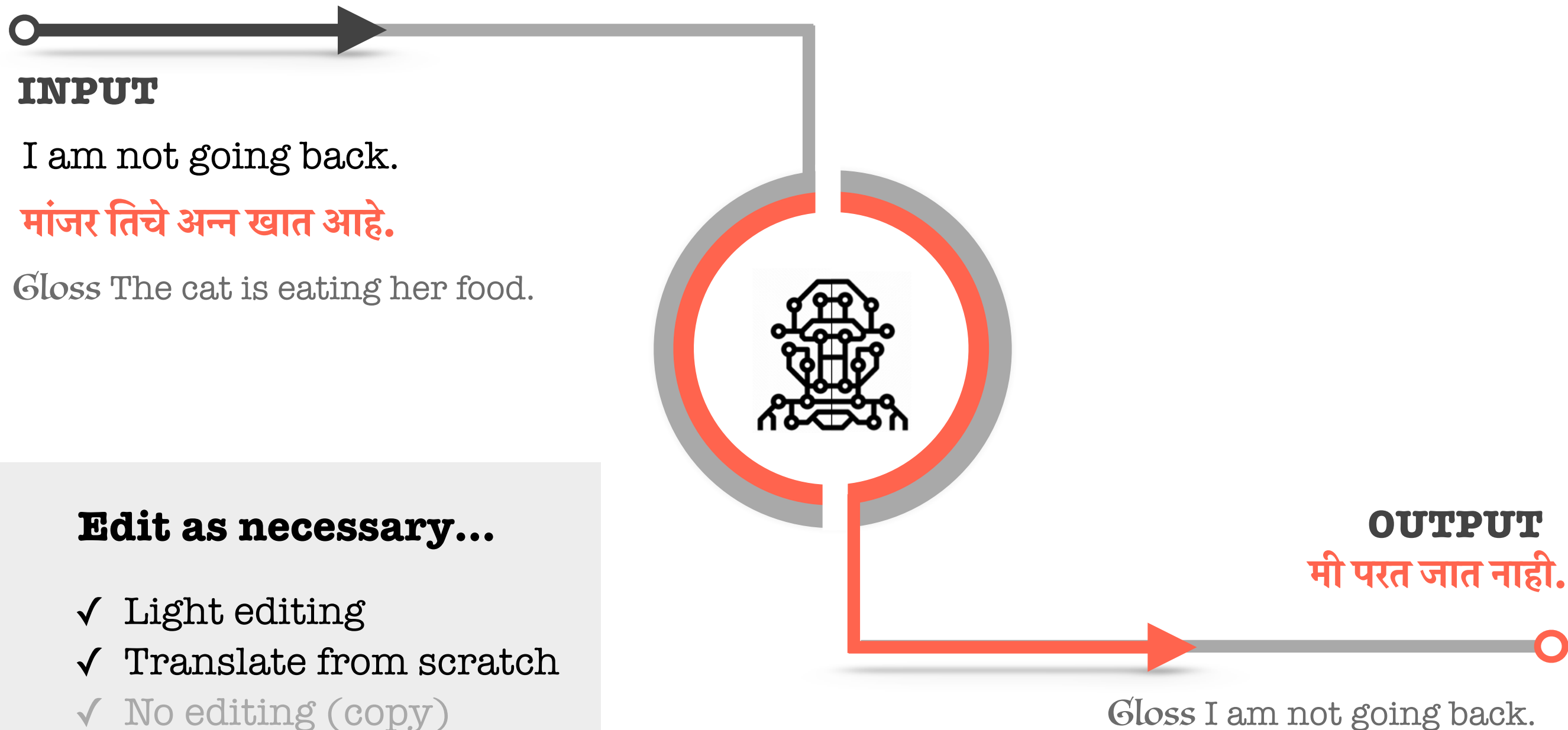
BitextEdit: Automatic Editing for Improved Bitext Quality



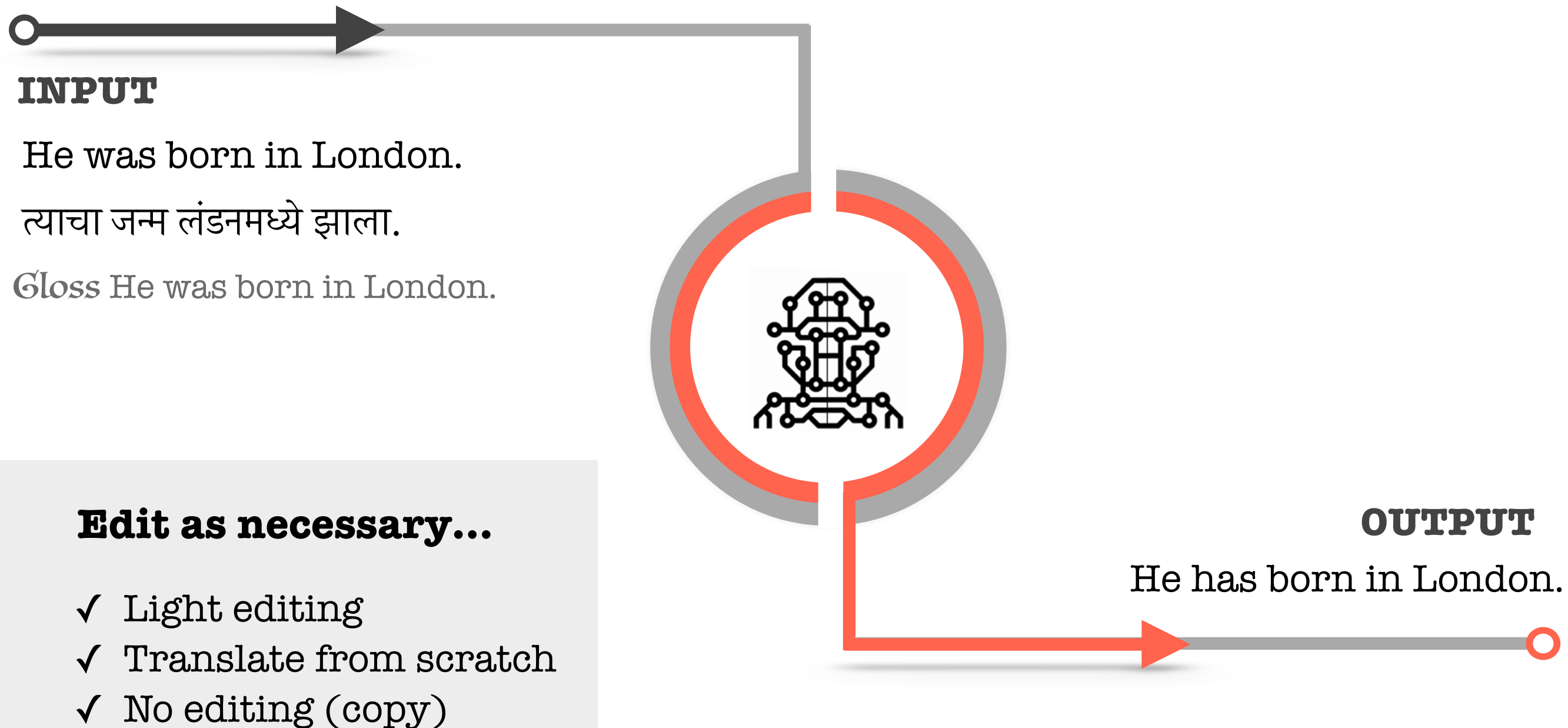
BitextEdit: Automatic Editing for Improved Bitext Quality



BitextEdit: Automatic Editing for Improved Bitext Quality

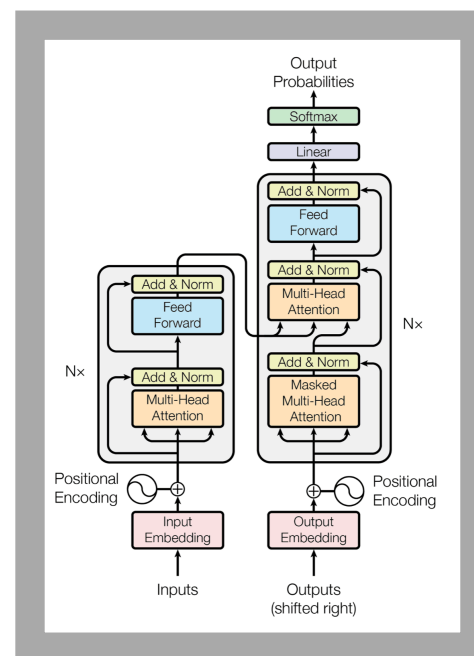


BitextEdit: Automatic Editing for Improved Bitext Quality



BitextEdit: Learn to Reconstruct Original from Noisy References & Translate

Sequence-to-Sequence Transformer



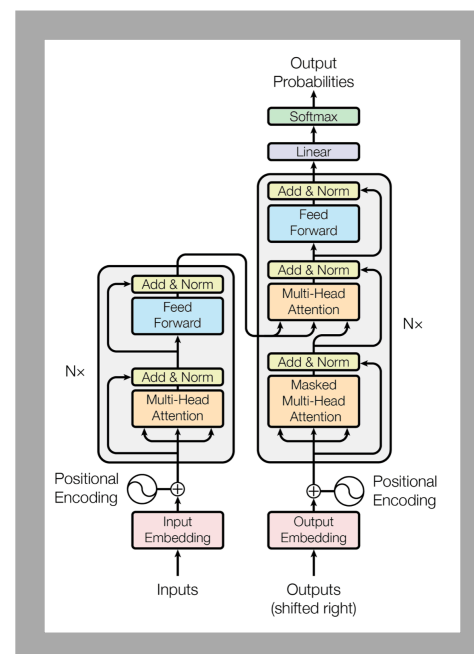
BitextEdit: Learn to Reconstruct Original from Noisy References & Translate

Sequence-to-Sequence
Transformer

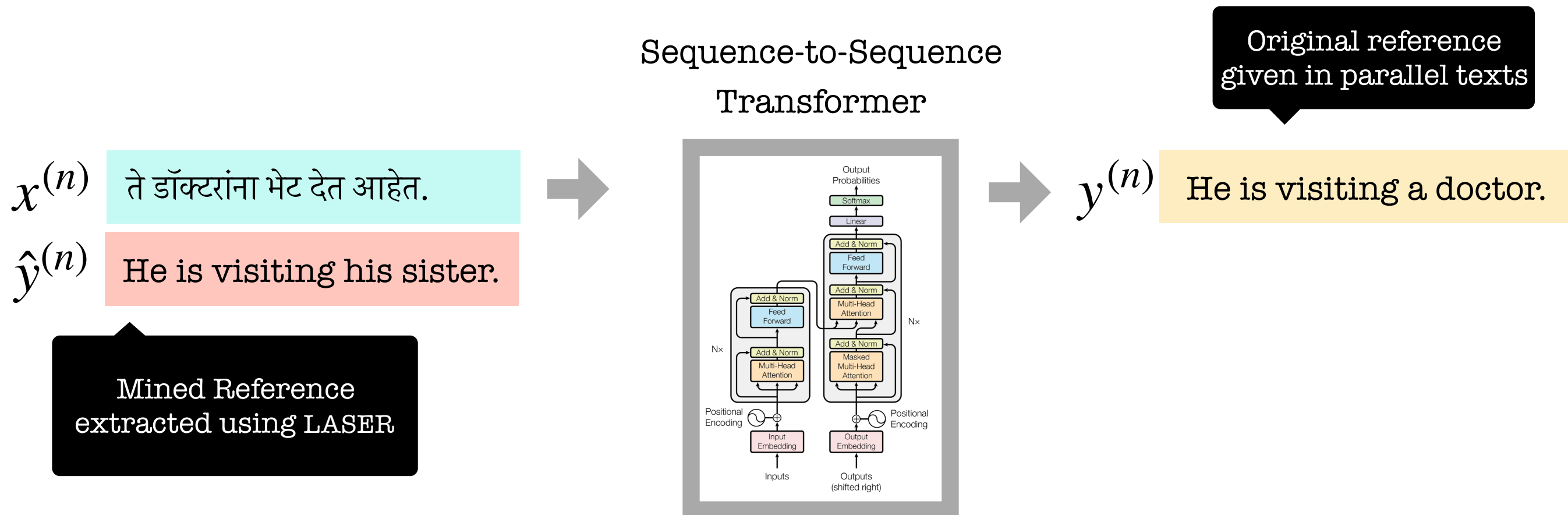
Original reference
given in parallel texts

$x^{(n)}$ ते डॉक्टरांना भेट देत आहेत.

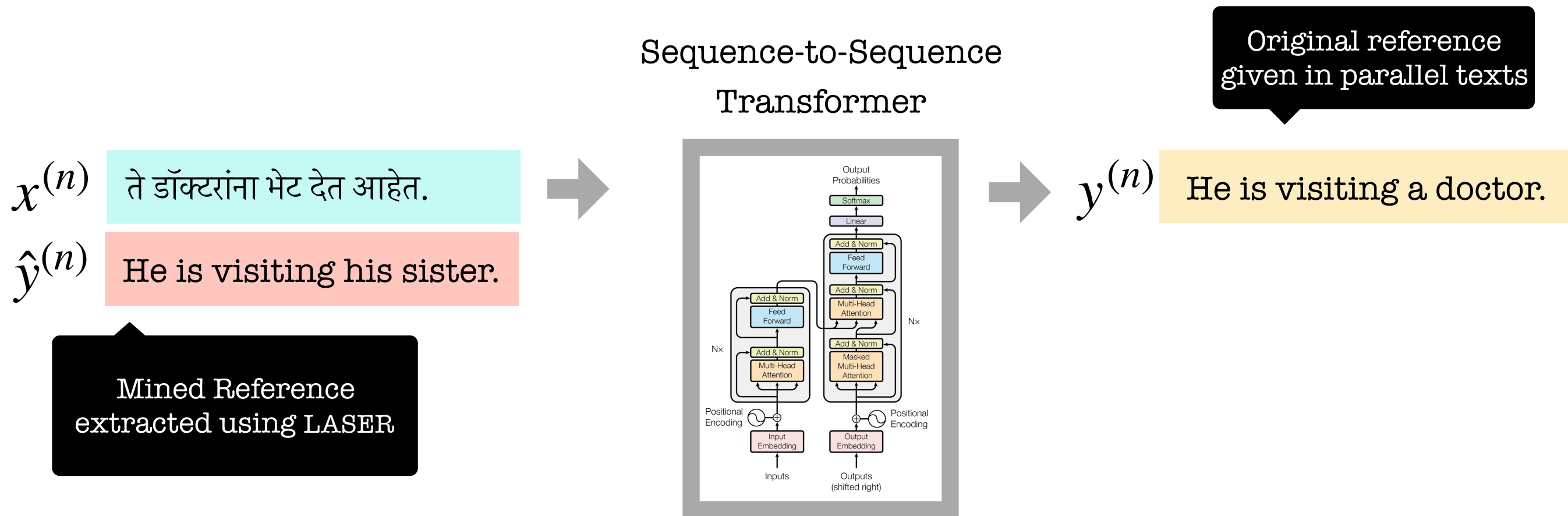
$y^{(n)}$ He is visiting a doctor.



BitextEdit: Learn to Reconstruct Original from **Noisy** References & Translate

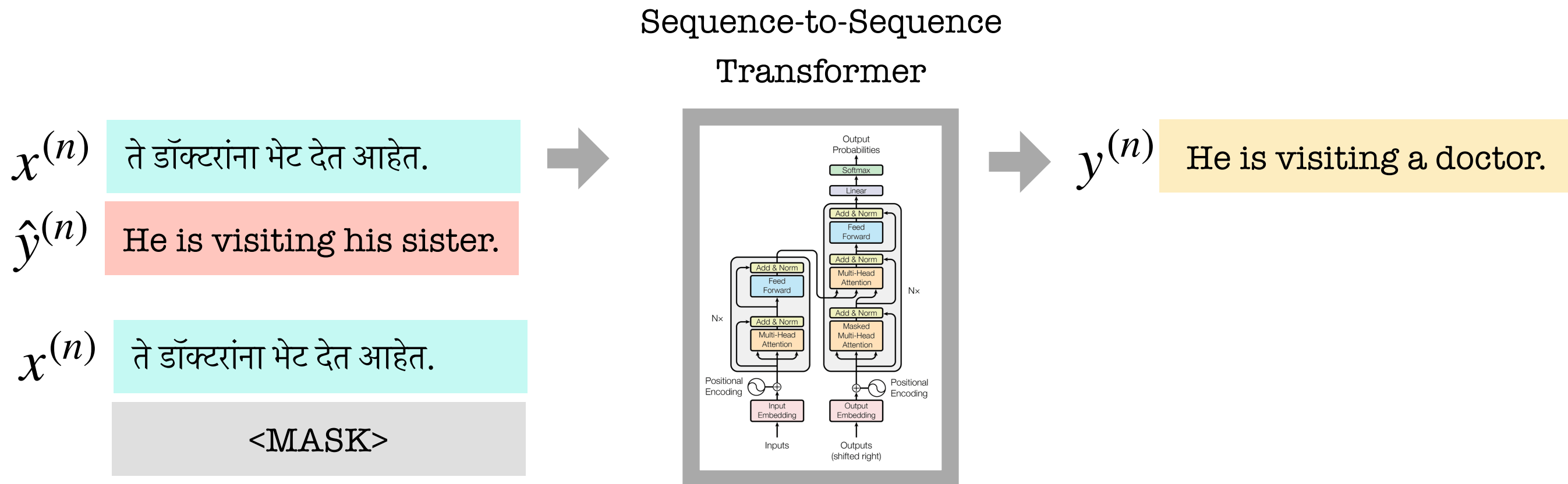


BitextEdit: Learn to Reconstruct Original from **Noisy** References & Translate



$$\log p\left([< e > y^{(n)}] | (x^{(n)}, \hat{y}^{(n)})\right) + \log p\left([< e > y^{(n)}] | (x^{(n)}, < \text{MASK} >)\right)$$

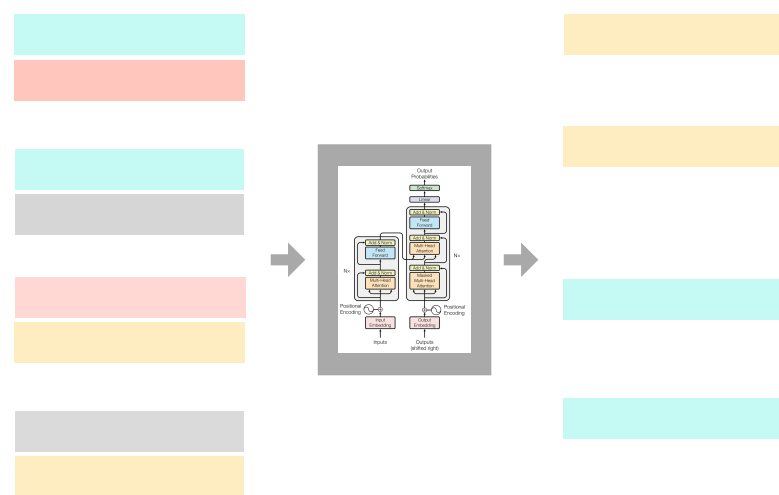
BitextEdit: Learn to Reconstruct Original from Noisy References & Translate



$$\log p\left([< e > y^{(n)}] | (x^{(n)}, \hat{y}^{(n)})\right) + \log p\left([< e > y^{(n)}] | (x^{(n)}, < \text{MASK} >)\right)$$

BitextEdit: Learn to Reconstruct Original from Noisy References & Translate

Bi-directional Training




$$\log p\left([< f > x^{(n)}] | (y^{(n)}, \hat{x}^{(n)})\right) + \log p\left([< f > x^{(n)}] | (y^{(n)}, < \text{MASK} >)\right) +$$

$$\log p\left([< e > y^{(n)}] | (x^{(n)}, \hat{y}^{(n)})\right) + \log p\left([< e > y^{(n)}] | (x^{(n)}, < \text{MASK} >)\right)$$

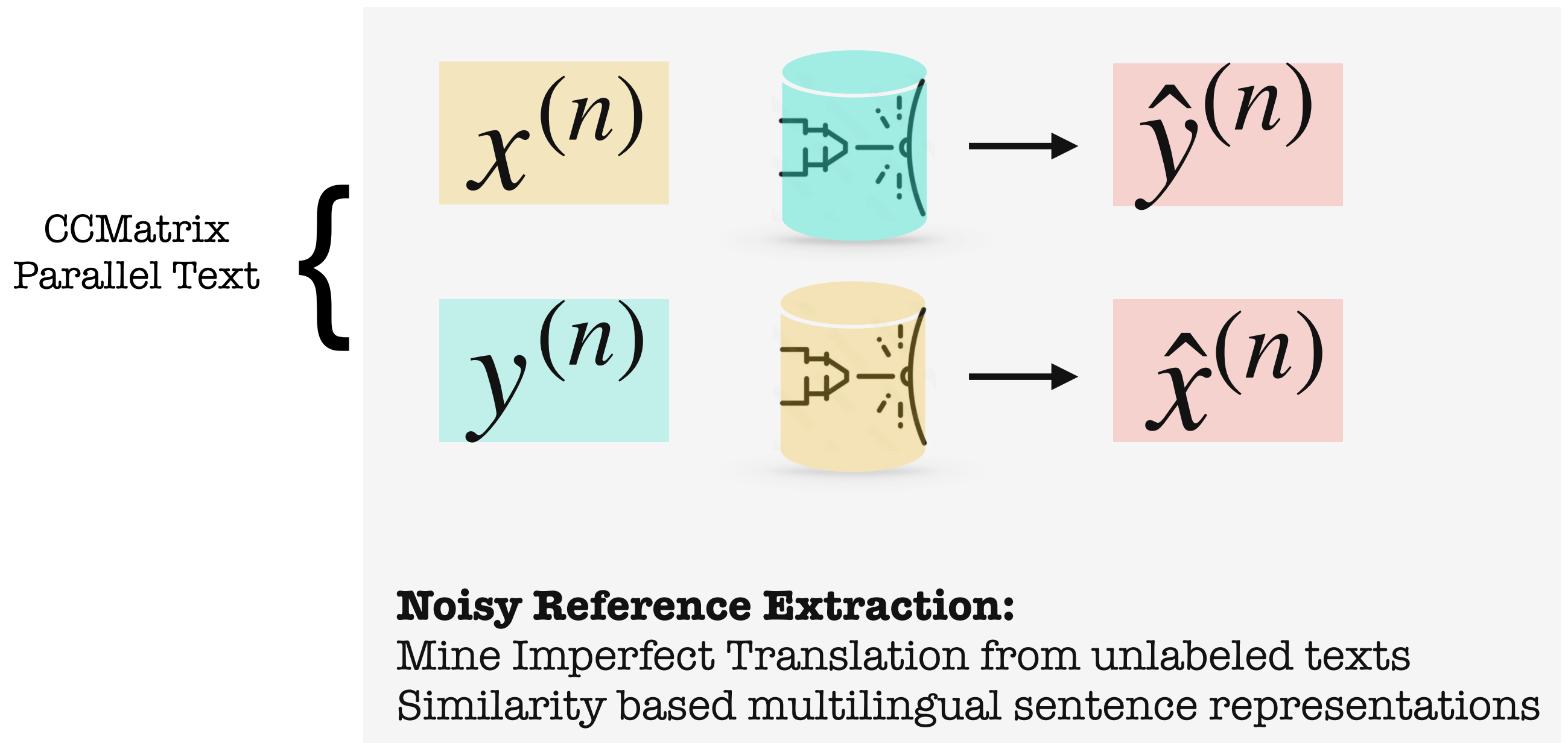
BitextEdit: Learn to Reconstruct Original from Noisy References & Translate

Where does this supervision $(x^{(n)}, y^{(n)}, \hat{x}^{(n)}, \hat{y}^{(n)})$ come from?



$$\begin{aligned} & \log p\left([\langle f \rangle x^{(n)}] \mid (y^{(n)}, \hat{x}^{(n)})\right) + \log p\left([\langle f \rangle x^{(n)}] \mid (y^{(n)}, \langle \text{MASK} \rangle)\right) + \\ & \log p\left([\langle e \rangle y^{(n)}] \mid (x^{(n)}, \hat{y}^{(n)})\right) + \log p\left([\langle e \rangle y^{(n)}] \mid (x^{(n)}, \langle \text{MASK} \rangle)\right) \end{aligned}$$

BitextEdit: Mine Potentially Imperfect Translations for each Text in Given Bitext

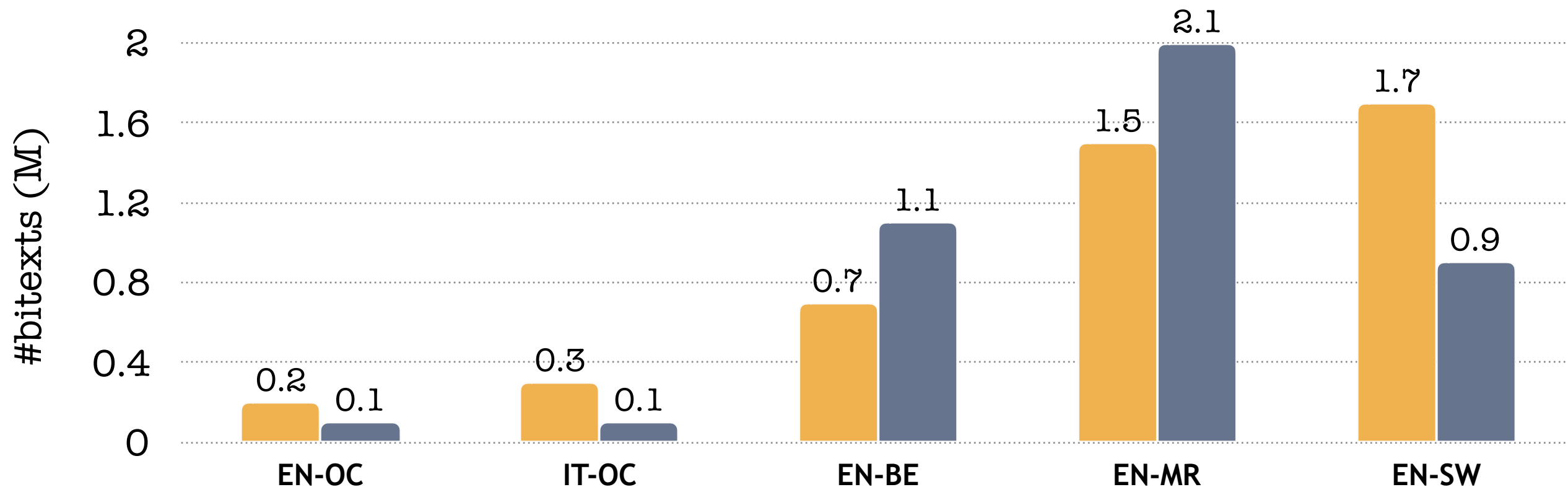


Experimental Settings

Bitexts that pass through
the noisy filtering step
(under LASER)

Bitexts that are filtered
based on the noisy filtering
step (under LASER)

PoolA
PoolB



Experimental Settings

Bitexts that pass through the noisy filtering step (under LASER)

Bitexts that are filtered based on the noisy filtering step (under LASER)

PoolA
PoolB

✓ **Training data:**

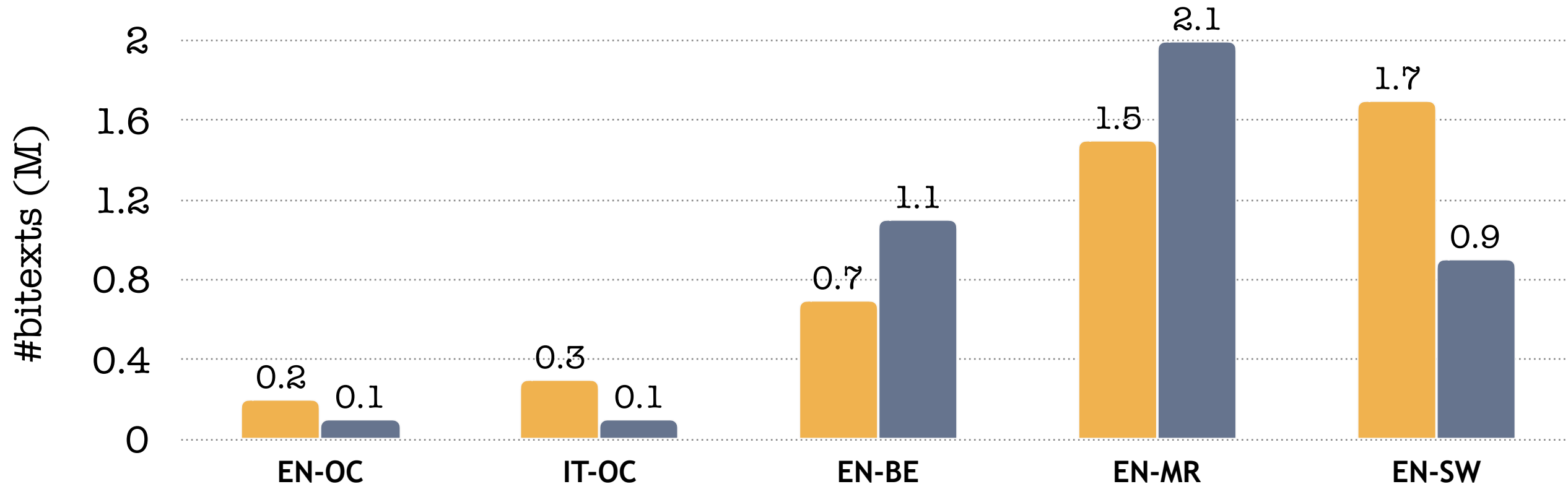
- ▶ CCMatrix Bitexts
- ▶ Revise Pool B

✓ **Model:**

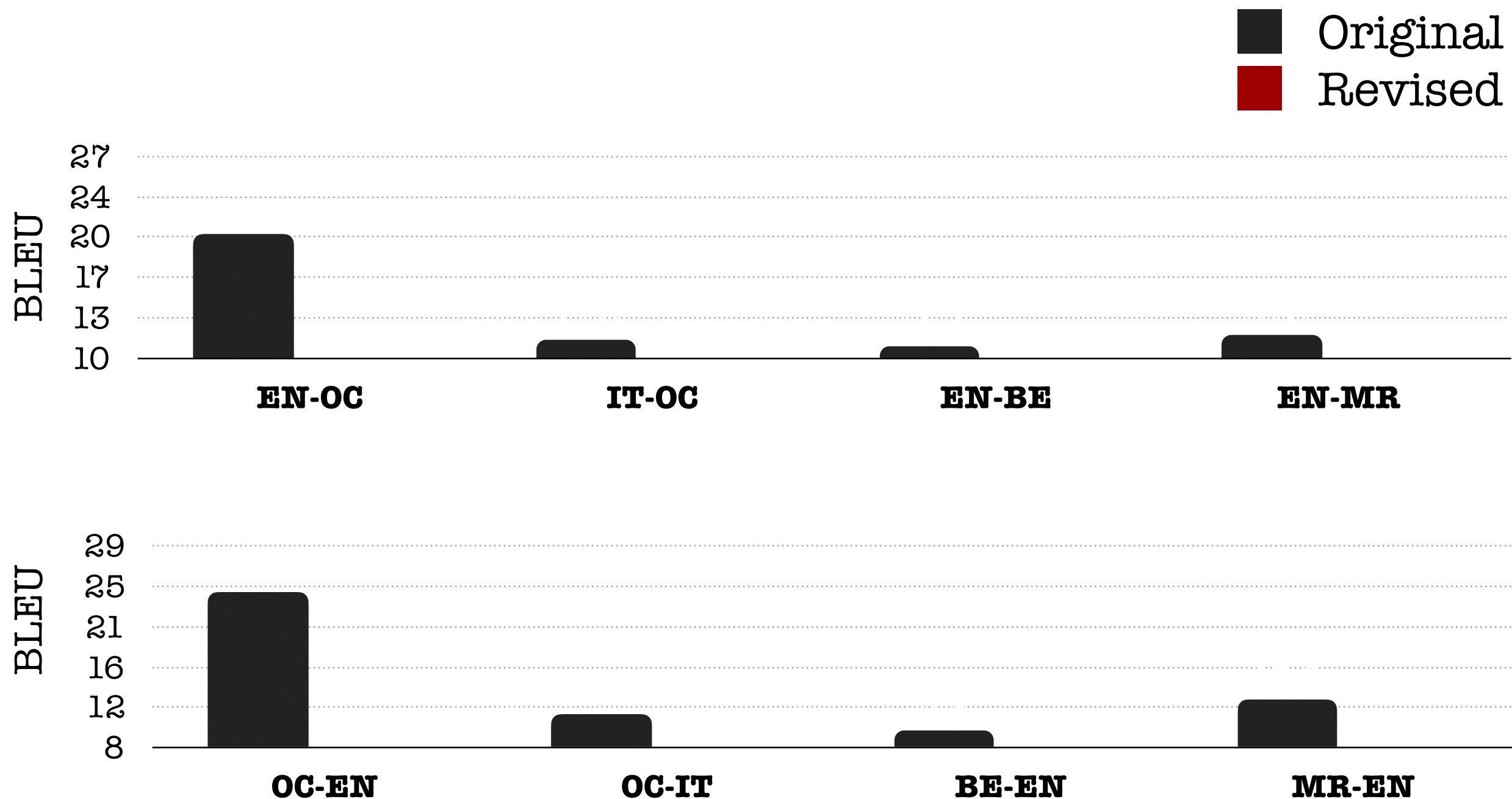
- ▶ Transformer NMT

✓ **Evaluation:**

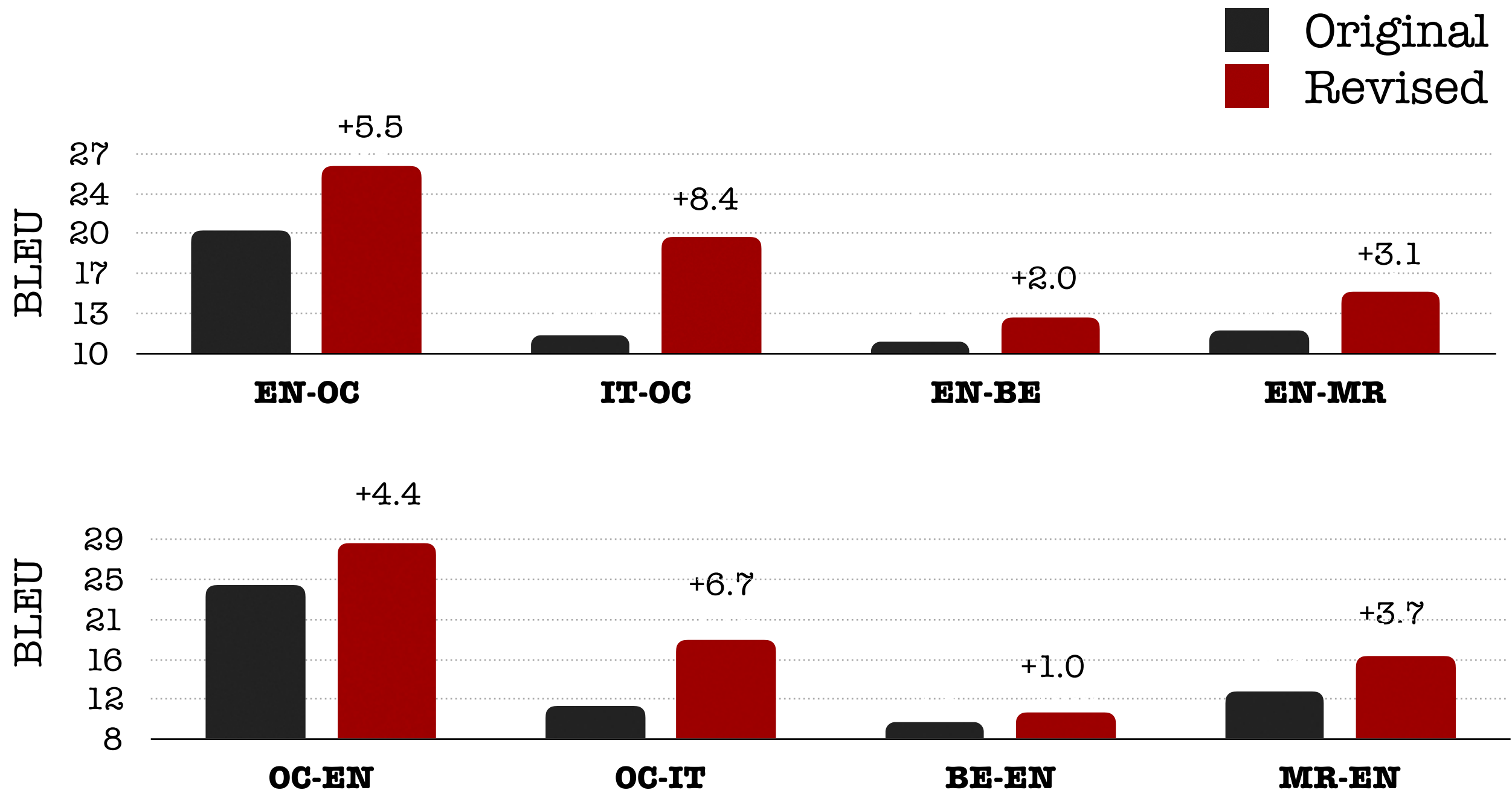
- ▶ BLEU (spm-BLEU)
- ▶ FLORES test set



Does the Revised Bitext Provide More Reliable Training Signal than the Original?



BitextEdit: Revised Bitext yields better Translation Quality than the Original



alternative approach for
bitext quality improvement

BitextEdit: Automatic Bitext Editing for Improved Low-Resource Machine Translation

synthetic supervision

alternative approach for
bitext quality improvement

BitextEdit: Automatic Bitext Editing for Improved Low-Resource Machine Translation

synthetic supervision

alternative approach for
bitext quality improvement

BitextEdit: Automatic Bitext Editing for Improved Low-Resource Machine Translation

translation quality
improvements on 10 tasks

synthetic supervision

alternative approach for
bitext quality improvement

BitextEdit: Automatic Bitext Editing for Improved Low-Resource Machine Translation

translation quality
improvements on 10 tasks

QUESTIONS?