

000 1 Appendix

001 In this appendix, we provide additional analysis,
002 results and examples.
003

004 1.1 Choice of languages and mined bitexts

005 We currently handle 90 languages. They were chosen
006 to cover several language families, frequent
007 high-resources languages as well as several low-
008 resource languages for which only a very limited
009 amounts of bitexts are publicly available. Mining
010 all possible $90 \times (90 - 1)/2 \approx 4000$ pairs
011 is computationally very challenging. In addition,
012 one may question if mining several “*unusual*” lan-
013 guage pairs would be useful, e.g. Icelandic-Urdu
014 or Gallican-Malay. It is quite unlikely that there is
015 interest in training a direct NMT systems for such
016 language pairs.
017

018 Therefore, we took the following approach. We
019 first organized all languages into twelve groups:
020

- 021 • Major Asian (4): Japanese, Korean,
022 Vietnamese, Chinese;
- 023 • Germanic (12): Afrikaans, Danish, Dutch,
024 German, English, Frisian, Icelandic, Luxem-
025 bourgish, Norwegian, Swedish, Yiddish;
- 026 • Romance (10): Asturian, Catalan, French,
027 Galician, Italian, Latin, Occitan, Portuguese,
028 Romanian, Spanish.
- 029 • Slavic (12): Belarusian, Bosnian, Bul-
030 garian, Croatian, Czech, Macedonian,
031 Polish, Russian, Serbian, Slovak, Slovenian,
032 Ukrainian;
- 033 • Other European (10): Albanian, Arme-
034 nian, Esperanto, Estonian, Finnish, Georgian,
035 Greek, Hungarian, Latvian, Lithuanian;
- 036 • Celtic/Irish (4): Breton, Irish, Scottish, Welsh;
- 037 • Turkic (5): Azerbaijani, Kazakh, Turkish,
038 Tatar, Uzbek;
- 039 • Middle East (3): Arabic, Farsi, Hebrew;
- 040 • Niger-Congo/Afro-Asiatic (10): Amharic,
041 Hausa, Igbo, Oromo, Somali, Swahili, Wolof,
042 Xhosa, Yoruba, Zulu;
- 043 • Indo-Aryan (9): Bengali, Hindi, Marathi,
044 Nepali, Oriya, Sinhala, Sindhi, Urdu, Tamil;

- 045 • Malayo-Polynesian (9): Cebuano, Ilocano,
046 Indonesian, Javanese, Malagasy, Malay,
047 Malayalam, Sundanese, Tagalog;
048

- 049 • Other Asian (2): Burmese, Khmer;
050

051 Most of these groups correspond to well estab-
052 lished linguistic language families, but we have
053 also performed some geographic groupings, in par-
054 ticular for small language families or isolated lan-
055 guages. We systematically mine all pairs within
056 one language family. For instance, we provide bi-
057 texts for all pairs of Indian languages. In addition,
058 we have identified major languages in each group
059 and use them as “*bridge languages*” (underlined in
060 the above list). We mine for all bitexts among these
061 27 bridge languages. The motivation for this bridge
062 language approach is to connect the languages of
063 the various groups, but still avoid mining the full
064 matrix.
065

066 We tried to support a large number of languages,
067 but we are aware that the underlying LASER em-
068 bedding is not very strong for all 90 languages,
069 for instance several languages of the Niger-Congo
070 family. Therefore, some of the mined bitexts may
071 contain wrong alignments or even texts from other
072 languages, despite careful filtering and two differ-
073 ent LID classifiers. Additional cleaning/filtering
074 may be needed. Nevertheless, we hope that these
075 bitexts are a useful resources to support research in
076 low-resource languages.
077

078 1.2 Example alignments

079 To illustrate the quality and richness of the mined
080 bitexts, we provide here some examples of ex-
081 tracted bitexts. We first searched for English sen-
082 tences which appear simultaneously in bitexts for
083 ten different languages (Arabic, German, French,
084 Indonesian, Japanese, Korean, Russian, Turkish
085 Vietnamese and Chinese). About ten thousand sen-
086 tences are such 11-way parallel. Table 1 gives four
087 examples. The first two examples are very generic
088 sentences which could appear on many Web pages.
089 This nicely showcases the potential of global min-
090 ing to find mutual translations in unrelated doc-
091 ments. The second example is rather long sentence
092 from some political document. Finally, we provide
093 an example from the medical domain. We also ob-
094 serve grammatically wrong sentences, e.g. “*Ein*
095 *Besuch in einem kranken Freund*”. This may indi-
096 cate that our approach can find parallel sentences
097 which were translated by (low quality) MT.
098

100				150
101				151
102				152
103				153
104				154
105				155
106	En	You should clean the refrigerator once a month.	Visiting a sick friend.	156
107	Ar	وأخوا ذكرى أنه يجب عليك تنظيف الثلاجة مرة واحدة في الشهر.	زرت صديقاً مريضاً.	157
108	De	Den Kühlschrank sollten Sie einmal im Monat saubermachen.	Ein Besuch in einem kranken Freund	158
109	Fr	Il est recommandé de nettoyer le réfrigérateur une fois par mois.	visite à un ami malade.	159
110	Id	Sebulan sekali kulkas harus dibersihkan.	Kunjungi teman yang sakit	159
111	Ja	1ヶ月に1回くらいは冷蔵庫の藏ざらえをしなきゃ。	病の友達を訪ねる	160
112	Ko	한 달에 한 번 정도는 냉장고 청소를 해주는 게 좋다.	아픈 친구를 보는 심정으로	161
113	Ru	Холодильник следует размораживать раз в месяц.	Посещение больного друга.	162
114	Tr	Buzdolabını boşaltarak ayda bir kez temizleyin.	Hasta bir dostu ziyaret etmek.	163
115	Vi	Vì vậy, mỗi tháng bạn nên vệ sinh tủ lạnh một lần.	Thăm người bạn THÀN bệnh	164
116	Zh	如果有必要，你可以一个月清理一次冰箱。	探望一个生病的朋友。	164

116	En	With the growing importance of world trade and the global community, business executives and legal professionals are expected to look beyond national jurisdictions and understand issues of international law and international commercial law.	166
117	Ar	مع تزايد أهمية التجارة العالمية والمجتمع العالمي، ومن المتوقع أن تنظر إلى أبعد السلطات القضائية الوطنية وفهم قضايا القانون الأوروبي والدولي المستشارين القانونيين.	167
118	De	Da Handel und Unternehmen immer globaler werden, wird erwartet, dass Rechtsberater über nationale Zuständigkeiten hinausblicken und Fragen des europäischen und internationalen Rechts verstehen.	168
119	Fr	Avec l'importance croissante du commerce mondial et la communauté mondiale, consultants juridiques devraient regarder au-delà des juridictions nationales et de comprendre les questions de droit européen et international.	170
120	Id	Dengan semakin pentingnya perdagangan dunia dan masyarakat global, konsultan hukum diharapkan untuk melihat melampaui yurisdiksi nasional dan memahami masalah hukum Eropa dan internasional.	171
121	Ja	法律コンサルタントは、貿易とビジネスがますますグローバル化するにつれて、国の管轄権を超えて、欧州および国際法の問題を理解することが期待されています。	172
122	Ko	무역 및 비즈니스가 전 세계적으로 증가함에 따라 법률 컨설턴트는 국가 관할권을 넘어서서 유럽 및 국제법 문제를 이해할 것으로 예상됩니다.	174
123	Ru	С ростом важности мировой торговли и мирового сообщества, юридические консультанты, как ожидается, искать за пределами национальной юрисдикции и понимания вопросов европейского и международного права.	175
124	Tr	Ticaret ve iş dünyası gittikçe küreselleştiğçe, hukuk müşavirlerinin ulusal yargılardan ötesine geçmesi ve Avrupa ve uluslararası hukuk konularını anlamaları beklenmektedir.	176
125	Vi	Với tầm quan trọng ngày càng tăng của thương mại thế giới và cộng đồng quốc tế, tư vấn pháp luật được dự kiến để nhìn xa hơn khu vực pháp lý quốc gia và hiểu các vấn đề của pháp luật châu Âu và quốc tế.	178
126	Zh	随着世界贸易和全球社会的重要性日益增加，法律顾问有望超越国家管辖和了解欧洲和国际法律的问题。	179

130	En	When we breathe quickly we also build up oxygen in our blood.	180
131	Ar	عندما نتنفس بسرعة نقوم ببناء الأكسجين في دمائنا.	182
132	De	Wenn wir schnell atmen, bauen wir auch Sauerstoff in unserem Blut auf.	183
133	Fr	Lorsque nous respirons rapidement, nous créons également de l'oxygène dans notre sang.	184
134	Id	Ketika kita bernapas dengan cepat, kita juga membangun oksigen dalam darah kita.	184
135	Ja	私たちが素早く呼吸すると、血液中に酸素も蓄積します。	185
136	Ko	우리가 빨리 숨을 쉬면 우리도 피 속에 산소를 축적합니다.	186
137	Ru	Когда мы дышим быстро, мы также накапливаем кислород в нашей крови.	187
138	Tr	Khi chúng ta thở nhanh, chúng ta cũng tích tụ oxy trong máu.	188
139	Vi	Çabucak nefes alduğumuzda, kanımızda da oksijen biriktiririz.	188
140	Zh	当我们快速呼吸时，我们的血液中也会积聚氧气。	189

Table 1: Examples of English sentences for which alignments in at least ten languages were found.

Model	Test set	de-en	en-de	en-ru	ru-en	zh-en	en-zh	de-fr	fr-de	
Ott et al. (2018)	Newstest2014	-	28.6	-	-	-	-	-	-	250 251
Fan et al. (2019)	Newstest2014	-	29.6	-	-	-	-	-	-	252 253
CCMatrix	Newstest2012	31.6	25.3	-	-	-	-	-	-	254
CCMatrix	Newstest2013	34.9	29.3	30.2	32.5	-	-	-	-	255
CCMatrix	Newstest2014	38.9	32.2	45.7	43.8	-	-	-	-	256
CCMatrix	Newstest2015	38.2	34.4	38.4	37.8	-	-	-	-	257
CCMatrix	Newstest2016	46.6	40.7	36.8	37.9	-	-	-	-	258
CCMatrix	Newstest2017	40.2	32.9	41.0	43.1	30.4	37.5	-	-	259
CCMatrix	Newstest2018	49.9	50.3	35.7	36.9	30.2	40.8	-	-	260
CCMatrix	Newstest2019	43.3	44.5	35.5	41.8	34.8	35.6	37.9	33.5	261
CCMatrix	Newstest2020	39.2	35.1	25.5	37.1	35.0	38.8	33.8	33.8	262

Table 2: Detokenized SacreBLEU scores of CCMATRIX models on all the available Newstest sets.

1.3 Additional WMT Results

We provide translation results as measured by SacreBLEU (Post, 2018) on all available WMT test sets in Table 2. For one of the most common translation evaluation benchmarks in the community, training on WMT16 en-de and evaluating on Newstest2014, we display the current state of the art results as well as the result of a well trained standard Transformer, to provide contrast against training on mined data only. On this benchmark, we find improvements of almost 2 BLEU points.

References

- Angela Fan, Edouard Grave, and Armand Joulin. 2019. Reducing transformer depth on demand with structured dropout.
- Myle Ott, Sergey Edunov, David Grangier, and Michael Auli. 2018. Scaling Neural Machine Translation. *arXiv:1806.00187*.
- Matt Post. 2018. A call for clarity in reporting bleu scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.