

Методы:

Название	Ссылка	Год	Комментарий	Этап (просмотрено/исс ледуется/не притрагивалась)
LLMs for Extremely Low-Resource Finno-Ugric Languages	<a href="https://openreview.net/attachment?id=KY3roODQ47&amp;name=pdf">https://openreview.net/attachment?id=KY3roODQ47&amp;name=pdf</a>	2024	<p>Статья полезна тем, что здесь описаны ресурсы, используемые для обучения (базис языков для pre-training: <b>русский, английский и латышский</b> (доля в данных по 12% каждый) и <b>финский с эстонским</b> по 32%). Если не будем предобучать из-за нехватки ресурсов, то стоит искать модели с таким бэкком.</p> <p>Вряд ли стоит рассматривать ллмки, у нас нет таких ресурсов, но вывод по bleu такой: ru-&gt;komi ~ 14.5</p>	Просмотрено
Machine Translation for Low-resource Finno-Ugric Languages	<a href="https://aclanthology.org/2023.nodalida-1.77.pdf">https://aclanthology.org/2023.nodalida-1.77.pdf</a>	2023	<p>Также поддерживает идею перевода монокорпуса и предлагает модели для “хорошей стартовой точки” для файнтюна (nllb и m2m).</p> <p>Использованы: M2M-100, 1.2 billion parameters (multi-lingual neural machine translation model);</p> <p>Для обучения также использовался Fairseq framework.</p> <p>Пример их finetune: <a href="https://github.com/TartuNLP/m2m-100-finetune">https://github.com/TartuNLP/m2m-100-finetune</a></p>	Просмотрено

			<p>Стоит также отметить, что они увеличивали вокабуляр и размер матрицы эмбедингов с помощью этих скриптов.</p> <p>Аккуратнее с библейскими данными, люди пишут, что переобучилось на них.</p>	
NEURAL MACHINE TRANSLATION FOR LOW RESOURCE LANGUAGES	<a href="https://arxiv.org/pdf/2304.07869">https://arxiv.org/pdf/2304.07869</a>	2023	<p>Используют перевод по словам для ускорения претренинга биязычной модели. Для претрейна биязычной модели используется masked language model (MLM) на моноязычных данных на обоих языках. Предлагают добавлять третий “язык” с переводом по словам, чтобы модель училась сопоставлять слова из обоих языков: “с целью MLM предсказать замаскированное английское слово, модель может учитывать как английские, так и иностранные слова в предложении “третьего языка”, и наоборот”.</p> <p>Идею для получения параллельных предложений с помощью меры жаккарда стоит попробовать на данных из Луима Серипос.</p> <p>Тоже использовался фреймворк Fairseq. Использовали самописный focal loss, описанный в <a href="https://aclanthology.org/2020.findings-emnlp.276.pdf">https://aclanthology.org/2020.findings-emnlp.276.pdf</a> (я также нашла реализацию <a href="https://github.com/vyraun/long-tailed/blob/main/fairseq/criterions/focal_loss.py">https://github.com/vyraun/long-tailed/blob/main/fairseq/criterions/focal_loss.py</a>) Тренировали mbart-cc25.</p>	Просмотрено

Machine Translation for Livonian: Catering to 20 Speakers	<a href="https://aclanthology.org/2022.acl-short.55.pdf">https://aclanthology.org/2022.acl-short.55.pdf</a>	2022	<p>Подводка к проекту OPUS с opensource кодом. Специализируются на лоу ресурс языки.</p> <p>Предлагаемая архитектура: 6 слоев энкодеров и декодеров, 8 attention heads на каждый слой, word embeddings и hidden layers размера 512, dropout на 0.3, максимальная длина предложения - 128 символов. Обучение проходило с помощью FairSeq тулы (pytorch) - <a href="https://github.com/facebookresearch/fairseq">https://github.com/facebookresearch/fairseq</a></p> <p>Их готовая модель для ливонского (как пример): <a href="https://huggingface.co/tartuNLP/liv4ever-mt">https://huggingface.co/tartuNLP/liv4ever-mt</a></p> <p>Генерация доп данных делается с помощью УЖЕ обученной модели, просто генеря лучшей моделью перевод монокорпуса.</p>	Просмотрено
Low-Resource Machine Translation Training Curriculum Fit for Low-Resource Languages	<a href="https://arxiv.org/pdf/2103.13272">https://arxiv.org/pdf/2103.13272</a>	2021	<p>Использовалась модель <a href="https://github.com/facebookresearch/XLM">https://github.com/facebookresearch/XLM</a>, потом был пре-трейн двуязычной LM на задаче MLM на монокорпусе. Также предлагается ввести третий язык, чтобы выровнять эмбединги английского и иностранного.</p> <p>Дальнейшая стадия включает в себя unsupervised пре-трейн NMT (энкодер и декодер - предобученная до этого модель) на бек-транслейшн монокорпуса.</p> <p>По доп данным также смотрят по Жаккарду. Также для оценки сопоставимых данных</p>	Просмотрено

			используется Ratio Margin-based Similarity Score.  Авторы отмечают, что используют 1гпу на 32Гб.	
Understanding Back-Translation at Scale	<a href="https://aclanthology.org/D18-1045.pdf">https://aclanthology.org/D18-1045.pdf</a>	2018	Создание синтетических данных для лоу ресурс	Не просмотрено
Exploring Diversity in Back Translation for Low-Resource Machine Translation	<a href="https://aclanthology.org/2022.deeplo-1.8.pdf">https://aclanthology.org/2022.deeplo-1.8.pdf</a>	2022	Тоже статья о работе с данными	Не просмотрено
Improving Low-Resource Neural Machine Translation with Filtered Pseudo-parallel Corpus	<a href="https://aclanthology.org/W17-5704.pdf">https://aclanthology.org/W17-5704.pdf</a>	2017	Берется предложение, выполняется перевод из таргета (моноязычные данные) в сурс-язык. получаем сурс_синтетику. переводим обратно из сурса_синт в таргет_синт. сравниваем таргет и таргет_синт по схожести. сортируем данные по метрике схожести, выбираем некоторое пороговое значение, все что ниже его - отсекаем. получаем синтетические данные для обучения. + используют бут	Просмотрено

Предлагаемые модели для безлайна:

Название	Ссылка	Комментарий	
smugri3-finno-ugric-nmt	<a href="https://huggingface.co/tartuNLP/s mugri3-finno-ugric-nmt">https://huggingface.co/tartuNLP/s mugri3-finno-ugric-nmt</a>	Весит 15Гб. Есть Манси, Ханты	
smugri3_14-finno-ugric-nmt	<a href="https://huggingface.co/tartuNLP/s mugri3_14-finno-ugric-nmt">https://huggingface.co/tartuNLP/s mugri3_14-finno-ugric-nmt</a>	Весит 9Гб. Есть Манси, Ханты	
tartuNLP/m2m-1_2B-finetune-finno-ugric-bt2	<a href="https://huggingface.co/tartuNLP/m2m-1_2B-finetune-finno-ugric-bt2">https://huggingface.co/tartuNLP/m2m-1_2B-finetune-finno-ugric-bt2</a>	Весит 5Гб. Файнтюн на финно-угрик языки.	
m2m-1_2B-finetune-finno-ugric-bt2	<a href="https://huggingface.co/tartuNLP/m2m-1_2B-finetune-finno-ugric-bt2">https://huggingface.co/tartuNLP/m2m-1_2B-finetune-finno-ugric-bt2</a>	Нет описания, авторы одни и те же с моделями выше, но весит чуть меньше 5Гб	
Opus-MT	<a href="https://github.com/Helsinki-NLP/Opus-MT-train?tab=readme-ov-file">https://github.com/Helsinki-NLP/Opus-MT-train?tab=readme-ov-file</a>		
facebook/nllb-200-distilled-1.3B	<a href="https://huggingface.co/facebook/nllb-200-distilled-1.3B/tree/main">https://huggingface.co/facebook/nllb-200-distilled-1.3B/tree/main</a>	Весит 5.5Гб. Есть примеры файнтюна в переводчик.	
facebook/m2m100_418M	<a href="https://huggingface.co/facebook/m2m100_418M/tree/main">https://huggingface.co/facebook/m2m100_418M/tree/main</a>	Весит 2Гб. Вроде эту модель использовали авторы кейса. Уже обучена как переводчик (на многие языки).	Большая вариативность глубины модели.
michaelfeil/ct2fast-m2m100_1.2B	<a href="https://huggingface.co/michaelfeil/ct2fast-m2m100_1.2B">https://huggingface.co/michaelfeil/ct2fast-m2m100_1.2B</a>	Несмотря на большое кол-во параметров, модель квантизирована и имеет сравнительно небольшой вес.	

tartuNLP/m2m100_418M_smugri	<a href="https://huggingface.co/tartuNLP/m2m100_418M_smugri/tree/main">https://huggingface.co/tartuNLP/m2m100_418M_smugri/tree/main</a>	Весит 2Гб. Предобучена на Финно-угорских языках.	
facebook/mbart-large-50	<a href="https://huggingface.co/facebook/mbart-large-50">https://huggingface.co/facebook/mbart-large-50</a>	Весит 2.5Гб. Много решений с использованием mBART.	Это предварительно обученная модель, которая в первую очередь предназначена для точной настройки в задачах перевода.
mBART	<a href="https://huggingface.co/facebook/mbart-large-cc25">https://huggingface.co/facebook/mbart-large-cc25</a>	Весит 2.5Гб. Много решений с использованием mBART.	Пример тюнинга mBART

Полезные гайды:

Название	Ссылки	Комментарий
Перевод	<a href="https://huggingface.co/docs/transformers/v4.17.0/en/tasks/translation">https://huggingface.co/docs/transformers/v4.17.0/en/tasks/translation</a> <a href="https://medium.com/@tskumar1320/how-to-fine-tune-pre-trained-language-translation-model-3e8a6aace9f">https://medium.com/@tskumar1320/how-to-fine-tune-pre-trained-language-translation-model-3e8a6aace9f</a>	Пример обучения из коробки, можно попробовать для быстрого безлайна