

Cross-lingual Multi-task Transfer for Zero-shot Task-oriented Dialog

Rob van der Goot¹, Marija Stepanovic¹, Alan Ramponi^{2,3}, Ibrahim Sharaf⁴,
Ahmet Üstün⁵, Aizhan Imankulova⁶, Siti Oryza Khairunnisa⁶, Mamoru Komachi⁶,
Barbara Plank¹

IT University of Copenhagen¹, University of Trento², Fondazione The Microsoft Research –
University of Trento Centre for Computational and Systems Biology (COSBI)³, Mawdoo3 AI⁴,
University of Groningen⁵, Tokyo Metropolitan University⁶
robv@itu.dk

1 Digital assistants for low-resource languages

Digital assistants are becoming an integral part of everyday life. However, commercial digital assistants are only available for a limited set of languages¹. Because of this, a vast amount of people can not use these devices in their native tongue. In this work, we focus on two core tasks within the digital assistant pipeline: intent classification and slot detection. Intent classification recovers the goal of the utterance, whereas slot detection identifies important properties regarding this goal (see an example in Figure 1). Besides introducing a novel cross-lingual dataset for these tasks, consisting of 11 languages, we evaluate a variety of models: 1) multilingually pretrained transformer-based models, 2) we supplement these models with auxiliary tasks to evaluate whether multi-task learning can be beneficial, and 3) annotation transfer with neural machine translation.

2 Data

For a long time, the Atis dataset (Hemphill et al., 1990) was used as main benchmark for task-oriented dialog tasks. It contains sentences from phone queries about flight information. However, since digital assistants have become more popular, focus has recently shifted to new domains (Schuster et al., 2019; Coucke et al., 2018). However, multilingual resources are scarce; the dataset from Schuster et al. (2019) contains data in Spanish and Thai (we refer to it as the FACEBOOK dataset). Very recently, Xu et al. (2020) have introduced a translation of the Atis dataset into 8 languages. In this work, we focus on the multilinguality of digital assistant queries, and translate more modern datasets (Schuster et al., 2019;

¹As of March 2020, between 8 to around 20 languages: <https://www.globalme.net/blog/language-support-voice-assistants-compared/>

Coucke et al., 2018) into 11 languages.

More specifically, we randomly sample 250 development and 150 test sentences from each of the datasets, and translate and annotate these to: Arabic, Danish, German, Indonesian, Italian, Kazakh, Dutch, Serbian, South-Tirolese (a German Southern-Bavarian dialect spoken in Northern Italy), Turkish and Chinese. We had one translator per language, who was also responsible for the annotation. All of our translators/annotators have a background in Natural Language Processing (NLP), are native speakers of the target language, and we had a balanced male/female distribution. During this annotation, we found some inconsistencies in the FACEBOOK data, and since the annotation guidelines were not publicly available, we wrote new annotation guidelines and re-annotated the English data as well.

Because we created new annotation guidelines for the FACEBOOK data, we let three annotators annotate the Dutch data. Their Fleiss Kappa score (Fleiss, 1971) was 0.924, indicating a near-perfect agreement. Common mistakes included annotation of question words, inclusion of locations in reminders, and the exact scope of the spans. We fixed the annotation after inspecting the disagreements, and updated the guidelines.

Figure 1 shows the annotation for one sentence in all our target languages from the Snips dataset (Coucke et al., 2018). We believe that this dataset contains a wide variety of language varieties, thereby providing a varied sample of different language phenomena.

3 Experiments

Baseline (BASE) We use MaChAmp (van der Goot et al., 2020), a NLP model focusing on multi-task learning implemented in AllenNLP (Gardner et al., 2018), as a baseline model. MaChAmp can exploit any HuggingFace embeddings (Wolf et al.,

AR	أود أن أرى مواعيد عرض فيلم Silly Movie 2.0 في دار السينما
DK	Jeg vil gerne se spilletiderne for Silly Movie 2.0 i biografen
DE	Ich würde gerne den Vorstellungsbeginn für Silly Movie 2.0 im Kino sehen
EN	I'd like to see the showtimes for Silly Movie 2.0 at the movie house
ID	Saya ingin melihat jam tayang untuk Silly Movie 2.0 di gedung bioskop
IT	Mi piacerebbe vedere gli orari degli spettacoli per Silly Movie 2.0 al cinema
KK	Мен Silly Movie 2.0 бағдарламасының кинотеатрда көрсетілім уақытын көргім келеді
NL	Ik wil graag de speeltijden van Silly Movie 2.0 in het filmhuis zien
SR	Želela bih da vidim raspored prikazivanja za Silly Movie 2.0 u bioskopu
DE-ST	I mecht es Programm fir Silly Movie 2.0 in Film Haus sechn
TR	Silly Movie 2.0 'in sinema salonundaki seanslarını görmek istiyorum
ZH	我想看 Silly Movie 2.0 在 影院 的放映

Figure 1: Examples of annotation for all languages in our dataset with intent: SearchScreeningEvent, and two slots: `movie_name` and `object_location_type`.

2019) as shared encoder, and uses a separate decoder for each task. In this work, we will use mBERT (Devlin et al., 2019) as encoder. We make use of both the SEQ and the CLASSIFICATION decoders in MaChAmp simultaneously to model both the slots and the intents.

Machine translation transfer (NMT.TRANSFER)

A commonly used approach is to translate the training data to the target language, map the annotation and then train a target-language model. We use the attention matrix to transfer the slot labels. This approach generally leads to a competitive performance (Schuster et al., 2019; Xu et al., 2020), but it is computationally costly to train machine translation models, and it is dependent on the size of the parallel data that is available. We use the Fairseq library (Ott et al., 2019) with default settings and use a combination of transcribed spoken parallel corpora, i.e., IWSLT 2016 Ted talks (Cettolo et al., 2016), Opensubtitles 2018 (Lison and Tiedemann, 2016)², and Tatoeba (Tiedemann, 2012).

Auxiliary tasks We experiment with a variety of auxiliary tasks to evaluate whether knowledge about these task can transfer to improve on our target tasks. For these auxiliary tasks, we use target language training data, and train simultaneously with the slots and intents decoders. The batches are mono-dataset, and we shuffle them on the batch level. The number of batches per epoch for both datasets is equal to the average number of batches per dataset. We use three different auxiliary tasks:

- Masked Language Modeling (AUX.MLM):

²<http://www.opensubtitles.org/>

	Facebook		Snips	
	slots	intents	slots	intents
BASE	64.1	65.0	81.0	94.9
NMT.TRANSFER	61.4	94.2	70.8	98.8
AUX.MLM	54.0	64.4	70.6	89.0
AUX.NMT	60.2	66.0	75.0	92.5
AUX.UD	65.8	64.9	71.0	78.8

Table 1: Average scores over all languages; f1-score (slots) and accuracy (intents).

we use the target language data from `nmt.transfer`, and mask words using the strategy from Devlin et al. (2019).

- Machine Translation (AUX.NMT): we use the data from NMT.TRANSFER, and use a recurrent neural network decoder with attention.
- Universal Dependencies (AUX.UD): We train on four tasks (UPOS, lemma, morph. tags, and dependency) of the Universal Dependencies (Nivre et al., 2020) simultaneously. We manually picked one treebank per language.

4 Results

Results are shown in Table 1. Initial results show that the baseline is hard to beat, and that the NMT.TRANSFER method is very strong for intents, and worse at slots because of mistakes in the alignment. Our proposed multi-task methods only outperform the baseline and the NMT.TRANSFER method for the slots in the FACEBOOK data. However, no tuning and/or analysis is done, leaving room for improvement.

References

- Mauro Cettolo, Niehues Jan, Stüker Sebastian, Luisa Bentivogli, Roldano Cattoni, and Marcello Federico. 2016. The IWSLT 2016 evaluation campaign. In *International Workshop on Spoken Language Translation*.
- Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, Maël Primet, and Joseph Dureau. 2018. [Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces](#). *arXiv:1805.10190 [cs]*. ArXiv: 1805.10190.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5).
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. [AllenNLP: A deep semantic natural language processing platform](#). In *Proceedings of Workshop for NLP Open Source Software (NLP-OSS)*, pages 1–6, Melbourne, Australia. Association for Computational Linguistics.
- Rob van der Goot, Ahmet Üstün, Alan Ramponi, and Barbara Plank. 2020. [Massive Choice, Ample tasks \(MaChAmp\): a toolkit for multi-task learning in NLP](#).
- Charles T Hemphill, John J Godfrey, and George R Doddington. 1990. The ATIS spoken language systems pilot corpus. In *Speech and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990*.
- Pierre Lison and Jörg Tiedemann. 2016. [OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 923–929, Portorož, Slovenia. European Language Resources Association (ELRA).
- Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. [Universal Dependencies v2: An evergrowing multilingual treebank collection](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4034–4043, Marseille, France. European Language Resources Association.
- 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 *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Sebastian Schuster, Sonal Gupta, Rushin Shah, and Mike Lewis. 2019. [Cross-lingual transfer learning for multilingual task oriented dialog](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3795–3805, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jörg Tiedemann. 2012. [Parallel data, tools and interfaces in OPUS](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Weijia Xu, Batool Haider, and Saab Mansour. 2020. [End-to-end slot alignment and recognition for cross-lingual NLU](#). *arXiv:2004.14353 [cs]*. ArXiv: 2004.14353.