

Few-Shot Learning in NLP

A Survey

Minh Duc Bui, mbui@mail.uni-mannheim.de

June 11, 2021

OUTLINE

Few-Shot Learning and Why its hard

Approaches

- Optimization-Based Meta-Learning Approaches

- Reformulate Tasks as language modelling problems

Few Shot Learning in Cross Lingual Setting

Summary & Discussion

Few-Shot Learning and Why its hard

WHAT IS FEW-SHOT LEARNING?

Supporting Set	
(A) capital_of	(1) <i>London</i> is the capital of <i>the U.K.</i> (2) <i>Washington</i> is the capital of <i>the U.S.A.</i>
(B) member_of	(1) <i>Newton</i> served as the president of <i>the Royal Society.</i> (2) <i>Leibniz</i> was a member of <i>the Prussian Academy of Sciences.</i>
Test Instance	
(A) or (B)	<i>Euler</i> was elected a foreign member of <i>the Royal Swedish Academy of Sciences.</i>

Figure 1: 2-Shot Relation Classification.

WHAT IS FEW-SHOT LEARNING?

Few-Shot Learning approaches **use prior knowledge** to generalize to new tasks containing only a **few samples** with supervised information.

Task T	Experience E		Performance P
	Supervised Information	Prior Knowledge	
Relation Classification	Few examples of Relations	Pre-learned language semantics	Acc.

TRANSFER LEARNING: SAMPLE EFFICIENCY

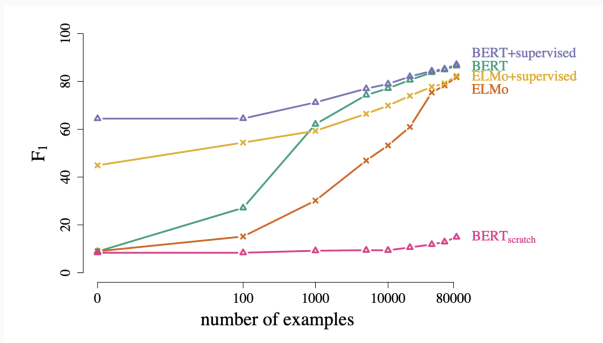


Figure 2: F1 scores on SQuAD as a function of the number of training examples (log scale). BERT+supervised denote BERT that is pretrained on other datasets and tasks.¹

¹Yogatama, D. et al. *Learning and Evaluating General Linguistic Intelligence*. 2019.

Approaches

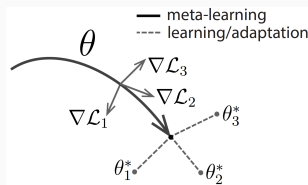
OPTIMIZATION-BASED APPROACHES

Goal:

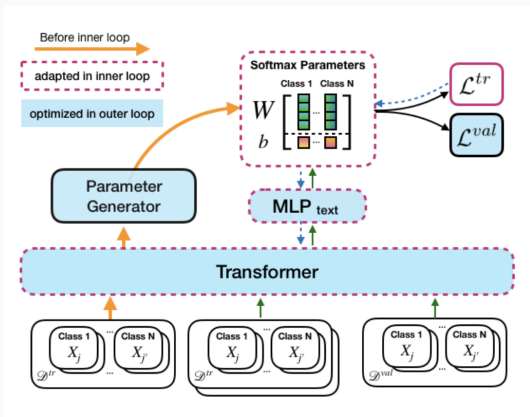
- **Learn a good set of parameter initialization** by using many tasks and treating **each task as a training example**

Training:

- Fine-tuning the model on a training set \mathcal{D}_i^{tr} of a selected training task, which **only consists of K examples**
- Use the task loss \mathcal{L}_i on \mathcal{D}_i^{test} to update our original not fine-tuned model parameters by computing the gradient



OPTIMIZATION-BASED: LEOPARD¹



Trained on 7 tasks and evaluated on 17 tasks in few-shot scenario

¹Bansal, T., Jha, R., and McCallum, A. "Learning to Few-Shot Learn Across Diverse Natural Language Classification Tasks". 2019.

OPTIMIZATION-BASED: DISCUSSION

- LEOPARD is first meta-learning approach that could generalize to test tasks, significantly different than training tasks (NLP)

N	k	BERT _{base}	MT-BERT _{softmax}	MT-BERT	Proto-BERT	LEOPARD
Overall Average	4	38.13	40.13	40.10	36.29	45.99
	8	36.99	45.89	44.25	39.15	50.86
	16	48.55	49.93	49.07	39.85	55.50

Figure 3: Few-shot generalization performance across tasks not seen during training.

META-LEARNING: DISCUSSION

Supporting Set	
(A) capital_of	
(B) member_of	
Test Instance	
(A) or (B)	<i>Euler</i> was elected a foreign member of <i>the Royal Swedish Academy of Sciences</i> .

Figure 4: 2-Shot Relation Classification. Can you do zero-shot learning?

HOW TO USE PRE-TRAINED LANGUAGE MODELS?

- Transformers are simply pre-trained on a language modeling objective in a semi-supervised fashion

Sentence:

The dog was chewing on a [MASK].

Mask 1 Predictions:

45.2% **bone**
30.1% **stick**
15.3% **toy**
9.4% **shoe**

- Reformulate Tasks as language modeling problems!**

REFORMULATE TASKS: FEW-SHOT WITH GTP-3¹

- Model is given a task description and K examples of context and completion, which they call model *priming*
- To make predictions, one final context is given, but the model has to fill in the completion

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

1	Translate English to French:	← task description
2	sea otter => loutre de mer	← examples
3	peppermint => menthe poivrée	
4	plush girafe => girafe peluche	
5	cheese =>	← prompt

¹ Brown, T. B. et al. *Language Models are Few-Shot Learners*. 2020.

REFORMULATE TASKS: GTP-3 RESULT

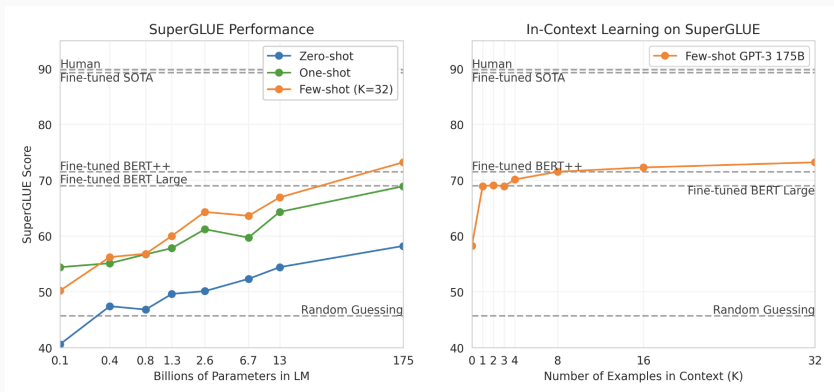


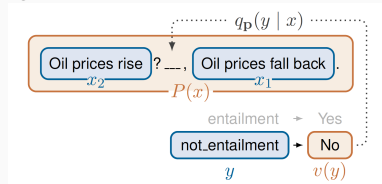
Figure 5: Performance on SuperGLUE increases with model size and number of examples in context.¹

¹ Brown, T. B. et al. Language Models are Few-Shot Learners. 2020.

REFORMULATE TASKS: OTHER APPROACHES

• iPET¹

- Reformulates tasks as cloze questions and uses ALBERT with regular gradient-based finetuning
- Uses knowledge distillation, which in turn need unlabeled data



	Model	Params (M)	BoolQ Acc.	CB Acc. / F1	COPA Acc.	RTE Acc.	WiC Acc.	WSC Acc.	MultiRC EM / F1a	ReCoRD Acc. / F1	Avg
test	GPT-3	175,000	76.4	75.6 / 52.0	92.0	69.0	49.4	80.1	30.5 / 75.4	90.2 / 91.1	71.8
	PET	223	79.1	87.2 / 60.2	90.8	67.2	50.7	88.4	36.4 / 76.6	85.4 / 85.9	74.0
	iPET	223	81.2	88.8 / 79.9	90.8	70.8	49.3	88.4	31.7 / 74.1	85.4 / 85.9	75.4
	SotA	11,000	91.2	93.9 / 96.8	94.8	92.5	76.9	93.8	88.1 / 63.3	94.1 / 93.4	89.3

Figure 6: Results on SuperGLUE for GPT-3 primed with 32 randomly selected examples and for iPET after training on 32 random examples.

¹ Schick, T. and Schütze, H. *It's Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners.* 2021.

Few Shot Learning in Cross Lingual Setting

PRE-TRAINING ON SIMILAR TASKS

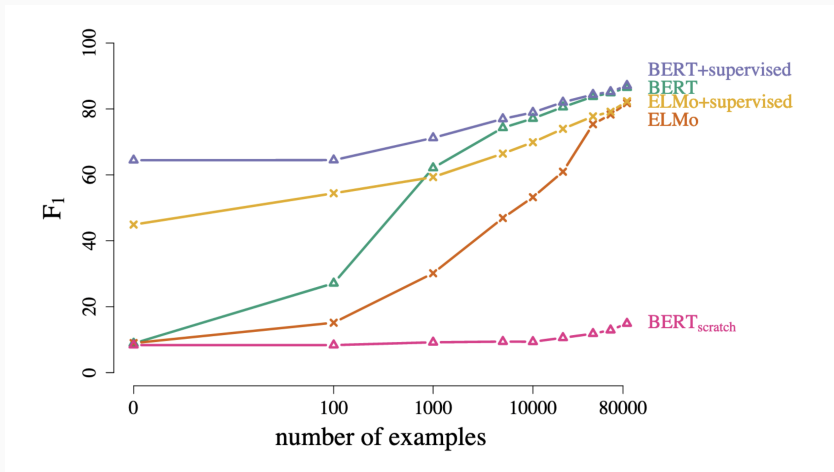


Figure 7: F1 scores on SQuAD as a function of the number of training

CROSS-LINGUAL SETTING

- Transfer the knowledge about the same task **from a high resource language to low resource language**
- Process:
 1. Couple a multilingual Transformer with task-specific classifier
 2. **Fine-tune model using task-specific supervised training data** from one **high resource language** (*source-adaption*)
 - If stop here: Zero-Shot Transfer
 3. **Continue fine-tuning** on K task-specific examples in the (low resource) target language (*target-adaption*).

CROSS-LINGUAL RESULTS: ZERO SHOT

Fine-tuning \ Eval	EN	DE	ES	IT
EN	96.82	89.40	85.91	91.60
DE	83.99	93.99	86.32	88.39
ES	81.64	88.87	96.71	93.71
IT	86.79	87.82	91.28	98.11

Figure 8: Zero Shot POS accuracy.¹

Task	Model	EN	ZH	TR	RU	AR	HI	EU	FI	HE	IT	JA	KO	SV
		Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ
DEP	B	92.3	-40.9	-41.2	-23.5	-47.9	-49.6	-42.0	-26.7	-29.7	-10.6	-55.4	-53.4	-12.5
POS	B	95.5	-33.6	-26.6	-9.5	-32.8	-33.9	-28.3	-14.6	-21.4	-6.0	-47.3	-37.3	-6.2
NER	B	92.3	-31.5	-6.5	-9.2	-29.2	-12.8	-8.5	-0.9	-9.2	-0.8	-51.1	-12.9	-1.9

Figure 9: Zero-shot cross-lingual transfer performance with mBERT (B).²

¹ Pires, T., Schlinger, E., and Garrette, D. How multilingual is Multilingual BERT?. 2019.

² Lauscher, A. et al. From Zero to Hero: On the Limitations of Zero-Shot Cross-Lingual Transfer with Multilingual Transformers. 2020.

CROSS-LINGUAL RESULTS: FEW SHOT

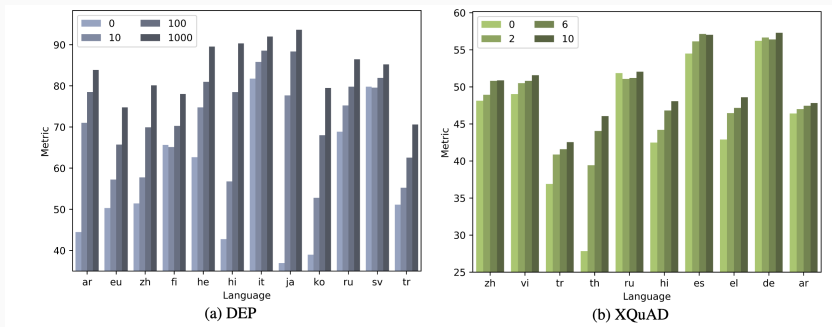


Figure 10: Results of the few-shot experiments with varying numbers of target-language examples k .¹

¹Lauscher, A. et al. From Zero to Hero: On the Limitations of Zero-Shot Cross-Lingual Transfer with Multilingual Transformers. 2020.

Summary & Discussion

BIG TRANSFORMERS, HIGH VARIANCE

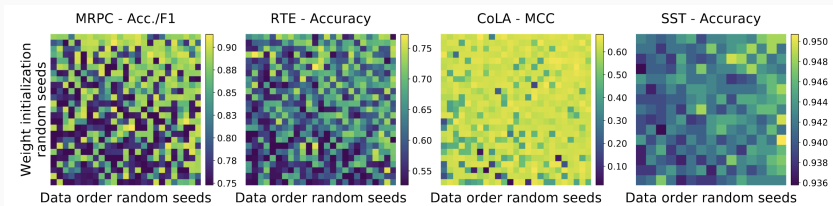


Figure 11: Visualization of validation performance, where each colored cell represents the performance of a training run¹.

¹Dodge, J. et al. *Fine-Tuning Pretrained Language Models: Weight Initializations, Data Orders, and Early Stopping*. 2020.

SUMMARY & DISCUSSION

- Discussed Optimization-based Meta-Learning
 - + Efficient use of few examples through optimal parameter initialization
 - Creating training tasks that enable finding a good initialization set to solve the target task is difficult.
 - Overfitting on task distribution

- Discussed reformulating to Language Modelling problems
 - + Achieve impressive results with small number of examples
 - Favor tasks, that can naturally be reformulated as "fill-in-the-blank"
 - Finding the right prompt is an art, needs a big enough validation set.¹
 - Restricted input size

¹ Gao, T., Fisch, A., and Chen, D. Making Pre-trained Language Models Better Few-shot Learners. 2020.

SUMMARY & DISCUSSION

- Discussed Few-Shot Cross-Lingual Transfer
 - + Finetuning on few examples, can significantly improve performances on distant languages - exactly where zero-shot CLT fail.
 - Inefficient way of finetuning.
 - Need labeled data in the low resource target language, which is typically hard to acquire.
- Few-Shot in General
 - + Enables to train without needing many training examples
 - + Advances to ultimate goal of NLP: General-purpose language understanding
 - Since most method rely on Transformers, they suffer high variance.
 - Essential to use the same set or average between multiple equal sets when comparing few-shot approaches, which is still lacking

Appendix

PRE-TRAINING ON LANGUAGE MODELLING OBJECTIVE

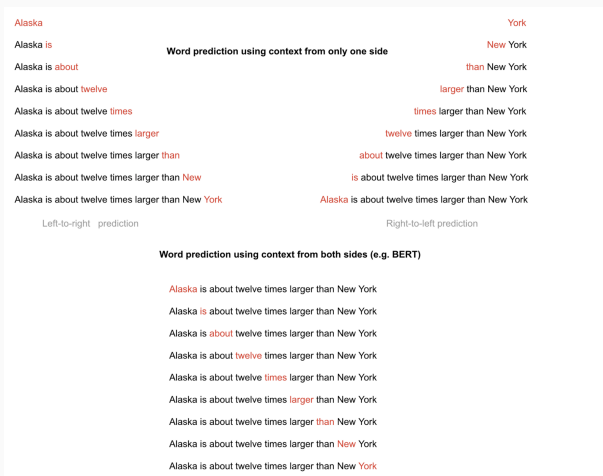


Figure 12: Pretraining on a language modelling objective.

PRIOR KNOWLEDGE: GENERAL PURPOSE LANGUAGE UNDERSTANDING

- Solving NLP tasks requires the model to learn about syntax, semantics, as well as certain facts about the world

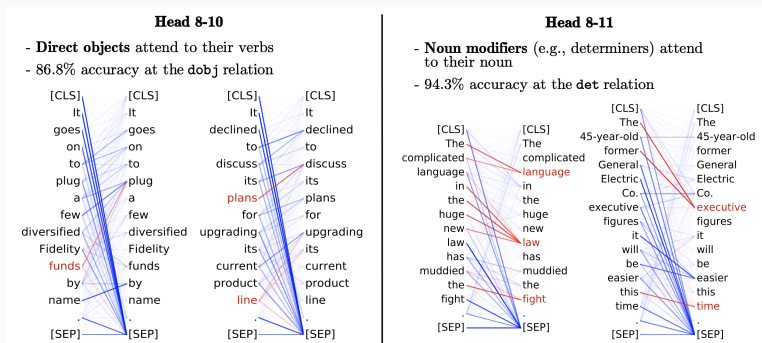


Figure 13: What Does BERT Look At? An Analysis of BERT's Attention [Clark et al. 2019]

META-LEARNING: DISCUSSION

- Challenge: Create training tasks enabling meta-learning algorithms to find a good initialization
 - Requires labeled data from many different tasks and additionally
- Major obstacle for meta-learning approaches is to solve diverse NLP few-shot tasks with one model
 - Suffers from overfitting to the training task-distribution (meta-overfitting)¹
 - Does not use any information of the underlying task

¹ Bansal, T., Jha, R., Munkhdalai, T., et al. "Self-Supervised Meta-Learning for Few-Shot Natural Language Classification Tasks". 2020.

REFORMULATE TASKS: DISCUSSION

- Achieve **impressive results** with only a small amount of examples
- Approaches **favor tasks, that can naturally be reformulated as "fill-in-the-blank" problems**, leaving room for future work
- LMs have **restricted input size**
 - Tasks that have long input sequences can not be properly solved
 - Future work: Use LM that allow long input sequences¹

¹ Beltagy, I., Peters, M. E., and Cohan, A. *Longformer: The Long-Document Transformer*. 2020.

BIG TRANSFORMERS, HIGH VARIANCE

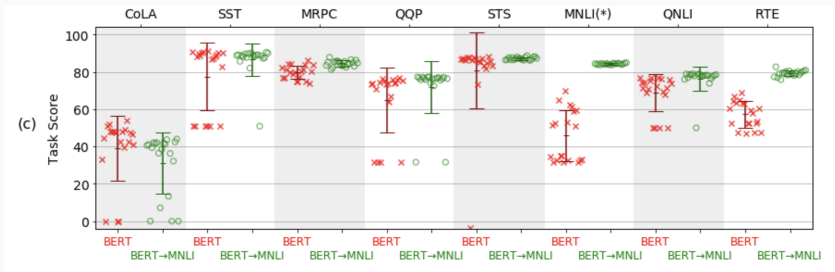


Figure 14: Distribution of task scores across 20 random restarts for BERT, and BERT with intermediary fine-tuning on MNLI. Fine-tuned on no more than 1k examples for each task.¹

¹ Phang, J., Févry, T., and Bowman, S. R. *Sentence Encoders on STILTs: Supplementary Training on Intermediate Labeled-data Tasks*. 2019.