

Few-Shot Training and Transfer in NLP

Anonymous EMNLP submission

Abstract

Natural Language Processing (NLP) tasks are data-hungry and when the situation arises, where data is scarce, NLP models often fail to carry out reliable generalizations. Humans can, however, generalize only by seeing a few labeled examples on a specific task. Motivated by this, the rise in popularity in techniques that can generalize to new tasks containing only a few samples, called Few-Shot Learning, was inevitable. This survey discusses pre-trained Language Models and Meta-Learning for Few-Shot Training and Transfer in NLP, critically assess their application and identifies future work. Furthermore, we study the application of Few-Shot approaches in a cross-lingual setting.

1 Introduction

Deep learning methods defined NLP in recent years, achieving impressive performance when sufficient amounts of labeled data are available. However, from a practical view, in many tasks, a large scale dataset is not available, e.g. low-resource languages, and annotating new labeled data labels is expensive and time-consuming (Fort, 2016), leaving us with only building more efficient algorithms as conventional deep learning methods fail in this low data regime (Yogatama et al., 2019). Humans, on the other hand, only need a few demonstrations to learn new language tasks. Motivated by this, Few-Shot Learning tries to solve all those issues by learning just from a few labeled samples.

1.1 Few-Shot Scenario

Few-Shot Learning (FSL) is the ability to learn tasks with limited examples. Most existing FSL problems are supervised learning problems, which is our focus in this survey. In an $(N\text{-way-})K\text{-shot}$ classification problem, we are only given K labeled examples per class, where the number of classes is N . $K\text{-shot}$ regression estimates a regression function given only K input-output example pairs.

To understand the challenges and approaches to Few-Shot Learning, we first analyze existing State-of-the-Art (SOTA) supervised approaches for NLP tasks.

1.2 Inducing Prior Knowledge

In a normal supervised setting, we would train our model on hundreds of thousands to millions of input-output pairs, which found success in numerous fields. However studies show that in NLP tasks, this paradigm of supervised learning does not generalize well outside the training data characteristics (Jia and Liang, 2017; Belinkov and Bisk, 2018), even when provided with enormous training data. The models are sensitive to noise, adversarial examples and are prone to overfitting. The reason is that language is complex and diverse and when conditions change, e.g. a new domain, the model is not able to adapt. Without any modification to the supervised approach, our Few-Shot Learning scenario will even amplify the poor generalization.

The most prominent way to help generalization is to induce an inductive bias by using transfer learning (Ruder, 2019), especially using pre-trained representations. In the last years, NLP saw the rise of pre-trained language representations for downstream tasks, achieving new SOTA on many NLP tasks. First, single-layer representations using word embedding vectors (Mikolov et al., 2013a) followed by contextualized word embeddings (Dai and Le, 2015; McCann et al., 2018a; Peters et al., 2018) were proposed, which were both simply fed into a task-specific architecture. With the rise of transformer language models (Vaswani et al., 2017), which enable direct finetuning of the whole architecture, there was no need for task-specific architectures anymore (Devlin et al., 2019). This was a breakthrough for NLP as many SOTA on NLP tasks were achieved by finetuning on task-specific samples using transformers, that are simply pre-trained on a language modeling objective

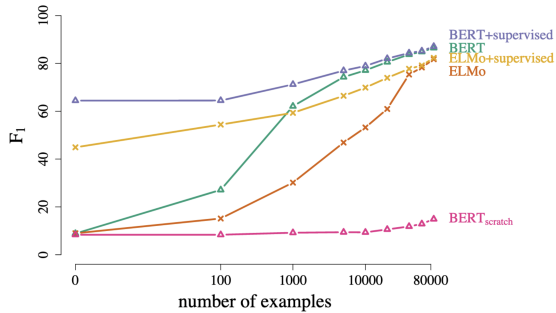


Figure 1: F_1 scores on SQuAD as a function of the number of training examples (log scale) (Yogatama et al., 2019). $BERT_{+supervised}$ and $ELMo_{+supervised}$ denote BERT and ELMo models that are pre-trained on other similar tasks, $BERT_{scratch}$ denotes a Transformer with a similar architecture to BERT that is trained from scratch.

in a semi-supervised fashion, inducing contextualized word embeddings. Clark et al. (2019) show that a pre-trained transformer model, like BERT, obtain knowledge about characteristics of the language, e.g. syntax and semantics as well as certain facts about the world, in short, have some general-purpose language understanding capacity, which can explain the generalization ability on a finetuned task.

1.3 Few-Shot Learning Challenges

One could naively apply the same strategy as in a normal supervised setting for our Few-Shot scenario: Finetune a pre-trained transformer model on the few labeled examples. Even though the pre-training helps the model to generalize on a new task, a sufficient amount of labeled data (Yogatama et al., 2019) is still needed in order to get reasonable results. The Figure 1 shows that in Few-Shot scenarios (i.e. < 1000 examples), all models lack far behind the fully trained variants, indicating *sample inefficiency* of the transformer model BERT. Additionally, the $BERT_{scratch}$ does not learn much without inducing an inductive bias via pre-training, showing the importance of the procedure. The Figure 1 also shows that, pre-training (sequentially) on similar tasks can help in a Few-Shot scenario, however, the sample inefficiency remains. Similar to this, Multi-Task Learning leverages information contained in multiple related tasks to help improve the generalization performance on all tasks (Zhang and Yang, 2021). Nonetheless, the method favors tasks with significantly more data, making it unsuitable for Few-Shot tasks. Another problem of big

transformers is that they suffer from high variance (Phang et al., 2019; Dodge et al., 2020). This is amplified in a Few-Shot scenario, where Language Models only finetune on a few samples (Zhang et al., 2021; Zhao et al., 2020). Changing the set of training examples can result in significant performance differences. Therefore, it is essential to use the same set or average between multiple equal sets when comparing Few-Shot approaches, making it hard to compare different approaches. (Zhang et al., 2021) provide alternative practices to reduce instability.

The question remains, how a transformer model can effectively leverage the few given examples without suffering from high variance. Section 2 describes a method that tries to exploit the induced bias of pre-trained language models explicitly through using the model directly by reformulating the task as a language model problem. Driven by the results of pre-training sequentially on similar tasks, see Figure 1, Section 3 will analyze Meta-Learning Approaches, which also "pre-train" on similar tasks to induce an inductive bias with the goal to use the model in a Few-Shot scenario. Section 4 will cover the use-case of K -Shot Cross-Lingual Transfer. Finally, Section 5 will conclude this survey.

2 Reformulating Tasks as Language Modelling Problems

As pre-trained language models possess some general language purpose understanding, the idea is to solve Few-Shot Learning tasks through directly using the obtained linguistic knowledge by reformulating tasks as language modeling problems and then predicting labels as "fill-in-the-blank" tasks, sharing the same format as pre-training LMs.

2.1 Approaches

Brown et al. (2020) introduces **GPT-3**, which essentially uses the same model as GTP-2 (Radford et al., 2019), but scales the data, training time, and model to 175 billion parameters. On the contrary to standard finetuning that condition on the task on the algorithmic level $p(output|input)$, the idea of GPT series is to condition the model on the selected task $p(output|input, task)$, by inducing the task into the text sequence with a task description. A reading comprehension training example could be formulated as (answer the question, document question, answer).

To create a training set, they scraped web pages, but with a focus on document quality. The hope is that task formulations occur naturally in the dataset. [Brown et al. \(2020\)](#) explore different settings for *learning within context*, which means that during inference, the model is given a *prompt*, which consists of a task description and K examples of context and completion, which they call model *priming*. Then to make predictions, one final context is given, but the model has to fill in the completion. One important note is, that the model does not do any weight updates during inference, even after seeing the K examples, leaving room for more optimization. Additionally, K is upper bounded by the context windows size ($n_{ctx} = 2048$), meaning that typically the window fits around 10 to 100 examples. With this strategy of using the pre-trained language model directly, GTP-3 shows impressive Few-Shot capabilities across diverse tasks, surpassing some strong finetuned models baselines, such as tasks in the SuperGLUE benchmarks ([Wang et al., 2020](#)) by only giving 32 labeled examples. However, for finding the right prompt, a hold out set is necessary, which then in return needs more examples. As we are in a Few-Shot scenario, this makes it difficult to obtain a sufficiently large hold out set. As GTP-3 naively concatenates the K randomly selected examples (as the model’s input size is bounded) with the input to create their in-context learning, the model does not make sure that the most informative demonstration are prioritized. However, prioritizing is important, since the number of usable demonstrations is bounded by the model’s input size. We will call this problem *in-context selection problem*. Furthermore, as GPT-3 uses an autoregressive language model, experiments do not include any bidirectional architectures, even though [Raffel et al. \(2020\)](#) indicate that (finetuned) models benefit from such bidirectionality to solve NLP tasks. Finally, as GPT-3 has 175 billion parameters, performing inference is expensive and makes it impracticable for many applications.

[Schick and Schütze \(2021b\)](#) introduce **iPET**, a task-agnostic method for Few-Shot Learning that can perform on par with the GTP-3 model on the SuperGLUE dataset using a 785 times smaller Language Model, making the approach more "greener" and practical. Instead of providing prompts, as in the GPT models, iPET uses pattern-exploiting training (PET) ([Schick and Schütze, 2021a](#)), which

reformulates tasks as cloze questions (no additional context samples provided) with regular gradient-based finetuning. Additionally, the model utilizes gradient steps after seeing the K examples. For that, PET requires a pattern-verbalizer pairs (PVPs) $\mathbf{p} = (P, v)$, which maps the input x of a task to a cloze question formulation. They call this a pattern P . Then for each possible output y of the task, PET maps it to a single token, representing its task-specific meaning in the pattern, called verbalizer v . Now, given a pre-trained masked language model, we only have to check the probabilities of the mapped output $v(y)$ being the correct token at the masked position. To generate good PVPs on a small development set of held out tasks, PET uses a combination of 3 PVPs per pattern for which a separate pre-trained MLM is first finetuned on the given (small) training set and then used to annotate unlabeled examples. Finally, the soft-labeled dataset is used to finetune a single sequence classifier, which is closely related to knowledge distillation ([Hinton et al., 2015](#)). However, PET only works when the answer is a single token. [Schick and Schütze \(2021b\)](#) proposes iPET, which modifies PET to handle more than just one token during predictions and refines the generation of PVPs by enabling them to learn from each other. [Schick and Schütze \(2021b\)](#) shows that iPET with ALBERT ([Lan et al., 2020](#)) as the underlying LM achieves similar results on the SuperGLUE dataset as GTP-3, given 32 examples. Additionally, iPET with ALBERT only uses 223 million parameters, which is a magnitude smaller than GTP-3. Even though iPET mitigates the problems of choosing a single cloze question formulation (pattern) by combining multiple formulations, it still requires engineering a set of suitable patterns. Furthermore, iPET requires additional unlabeled data, which it uses in the knowledge distillation stage. This can be hard to acquire, where samples are pairs of text with a label, constructed to test a model’s natural language understanding abilities (e.g. SuperGLUE). ([Tam et al., 2021](#)) propose **ADAPET**, which uses no unlabeled data by providing more supervision by modifying PET’s objective. ADAPET outperforms iPET on SuperGLUE without any unlabeled data.

GTP and PET models use prompt-based (pattern-based) prediction, but finding the right prompt/pattern is an art. [Gao et al. \(2020\)](#) proposes **LM-BFF**, which alleviates this problem by

generating the prompt automatically given a few examples, outperforming or matching manually selected prompts. Gao et al. (2020) first finds a label word mapping given a template (pattern) and then generates a diverse set of templates from the fixed set of label words by using the T5 model (Rafael et al., 2020). Even though Gao et al. (2020) propose a way to automatically find prompts, it still needs an "initial" template (pattern) or label words, inducing a bias that could potentially restrict the search space to a suboptimal one. Contrary to iPET, LM-BFF uses demonstrations for each input by concatenating them for additional context. However, to combat the *in-context selection problem* (see GTP-3), LM-BFF randomly selects a single example from *each* class for each input iterative at a time to create multiple, minimal demonstration sets, making it more efficient for Few-Shot tasks than GTP-3. As the underlying Language model, they use RoBERTa large model, which is again a magnitude smaller than GTP-3, with $K = 16$ examples and then use prompt-based finetuning with demonstrations. Notice that the finetuning process is different than iPET, which does not use any demonstrations, and GTP-3's in-context learning, which simply concatenates the input with demonstrations randomly drawn from the training set with no finetuning. Gao et al. (2020) evaluates on 8 tasks from the GLUE benchmark (Wang et al., 2019), SNLI (Bowman et al., 2015), and sentence classification tasks. Gao et al. (2020) show that their method of prompt-based finetuning outperforms standard finetuning (on $K = 16$ examples), except for the CoLA task and outperforms the GTP-3-style in-context learning. They also show that using demonstrations in context performs better than without any demonstrations in the context.

2.2 Discussion: Reformulating Tasks for NLP Few-Shot Tasks

Even though the models presented here, achieve impressive results with only a small amount of examples, it is still lacking quite far behind SOTA models that finetune on big datasets with thousands of examples. These approaches also favor tasks, that can naturally be reformulated as "fill-in-the-blank" problems, such as sentiment classification (e.g. positive class: "A fun ride. All in all **great**."), leaving room for future work. Additionally, methods require manual work to find a good reformula-

tion for tasks. This problem is amplified in practical situations, where we want to deploy such systems since we need domain and model expertise to find an optimal reformulation by hand for unseen tasks. Even though Gao et al. (2020) try to mitigate this problem by automatically find reformulations, LM-BFF still needs an initial reformulation. Additionally, Language Models have a restricted input size. Tasks that have too long input sequences can not be properly solved. Future work could investigate using Language Models that allow such long input sequences, e.g. Longformer (Beltagy et al., 2020). Furthermore, these approaches finetune the downstream tasks in isolation, not utilizing any information from similar tasks.

3 Meta-Learning Algorithms

Additionally to the general-purpose language understanding properties of pre-trained language models, Meta-Learning algorithms try to induce another inductive bias, which allows the model to quickly adapt after only seeing a few examples. In comparison to the methods described in Section 1.2, Meta-Learning explicitly take the Few-Shot scenario into account and utilize information from similar tasks. This is achieved by collecting many training tasks, where each training task consists of a training dataset $\mathcal{D}_i^{tr} = \{(\mathbf{x}_{tr}, \mathbf{y}_{tr})\}$, called support set S , and a test set $\mathcal{D}_i^{val} = \{(\mathbf{x}_{val}, \mathbf{y}_{val})\}$. The idea is to then pre-train on them such that the final model can generalize to new tasks rapidly, which allows us to perform Few-Shot tasks.

We will discuss 2 popular forms of Meta-learning for NLP tasks (Yin et al., 2020a): Metric-based and Optimisation-based learning.

3.1 Metric-based Meta-Learning

The idea in metric-based Meta-Learning is to learn a representation space through the training tasks, which enables us to classify test instances correctly by just comparing them to the K labeled examples in this representation space.

Vinyals et al. (2017) proposes **Matching Networks**, which use two different embedding functions, one for the training examples and one for the test examples. The representation of one example can change, depending on the given support set \mathcal{D}_i^{tr} for the task \mathcal{T}_i . For a test example \hat{x} , given its support set S , we choose the class with the highest aggregated similarity between class examples in the support set and the test instance by calculat-

ing the cosine similarity in the embedding space. Vinyals et al. (2017) evaluated Matching Networks on Few-Shot language modeling. Even though the approach found many applications in image classification, it has not yet found any impressive results in NLP tasks. One of the reasons is that matching networks do not finetune on the support set during inference, making it hard to find a good general embedding space that would work for many NLP tasks since text is quite diverse. If Matching Networks choose to finetune, it suffers from overfitting issues (Vinyals et al., 2017), not gaining much in performance.

To enable finetuning during inference and not suffer from overfitting, Snell et al. (2017) propose **Prototypical Networks**, which induce a simple bias: There exists an embedding space where points that belong to one class, cluster around a single prototype representation. For that, they learn a non-linear mapping of the input into an embedding space using a neural network and calculate the class’ prototype as the mean of its support set in the embedding space. Finally, we can classify a new instance by finding the nearest class prototype. In comparison to Matching Networks, they do not compare instances to each other but use the prototypes (class representation) calculated from the support set. Therefore, only in a Few-Shot scenario, the approaches differ. Additionally, Prototypical Networks use Euclidean distance which outperforms the proposed cosine similarity of Matching Networks (Snell et al., 2017). Prototypical networks were first originally suggested for images in computer vision problems, however, the method was also applied to NLP tasks. Most applications use pre-trained word embeddings and instead of averaging to calculate the prototype class, they use more sophisticated methods, such as attention-based prototypes, reaching new SOTA on some benchmarks (Han et al., 2018; Gao et al., 2019; Hui et al., 2020; Deng et al., 2020) and also finding applications in domain transfer (Bansal et al., 2019). However, as the metric plays an important role in gaining performance (Snell et al., 2017), Sung et al. (2018) introduces a learnable metric instead of a fixed metric, calling it **Relation Networks**. Yu et al. (2018) try to solve diverse Few-Shot text classification by extending Prototype Networks with clustering similar training tasks, learning one metric for each, and then automatically determining the best weighted combination of those metrics for

a newly seen Few-Shot task.

Matching, Prototypical and Relation Networks in NLP are mostly restricted to test tasks that are very similar to the training tasks, e.g. doing domain transfer. When we have diverse NLP tasks, finding an appropriate metric space becomes much harder. Yu et al. (2018) try to solve diverse Few-Shot text classification by extending Prototype Networks with clustering similar training tasks, learning one metric for each, and then automatically determining the best weighted combination of those metrics for a newly seen Few-Shot task. They show significant gains on Few-Shot sentiment classification and dialog intent classification tasks, indicating that clustering related tasks to handle diverse Few-Shot NLP tasks, might be a good research direction to improve metric-based or even optimization-based Meta-learning approaches for Few-Shot NLP tasks. A closely related method to metric-based approaches is supervised contrastive learning (Gunel et al., 2021) as both rely on capturing the similarity between examples in one class and contrasting them with examples in other classes. Instead of the usual Cross-Entropy Loss of Language Models, which is prone to high variance, Gunel et al. (2021) propose a loss function, consisting of cross-entropy and their supervised contrastive learning (SCL) term that pushes examples from the same class closer and the examples from different classes further apart. Gunel et al. (2021) obtain significant improvements over a strong RoBERTa-Large baseline on multiple datasets of the GLUE benchmark in few-shot learning settings. This method is closely related to metric-based approaches as both rely on capturing the similarity between examples in one class and contrasting them with examples in other classes.

3.2 Optimization-based Meta-Learning

In contrast to metric-based Meta-Learning, optimization-based Meta-Learning approaches try to learn a *good set of parameter initialization*, such that the model can quickly converge to a minimum in just a few gradient descent steps.

Finn et al. (2017) proposed the first optimization-based model, called **Model-Agnostic Meta-Learning (MAML)**. Let θ denote the parameter initialization of the model, ϕ_i the finetuned model parameters and \mathcal{L}_i the loss function of each task \mathcal{T}_i . The idea is to first sample a (or batch of) task \mathcal{T}_i with the corresponding (disjoint)

465 datasets $\mathcal{D}_i^{tr}, \mathcal{D}_i^{val}$. To train the model on \mathcal{D}_i^{tr} ,
 466 gradient-finetune with respect to the loss function
 467 \mathcal{L}_i to obtain ϕ_i . We then can update the original
 468 initial model parameters θ using the "test" loss
 469 $\mathcal{L}_i(\phi_i, \mathcal{D}_i^{val})$ across sampled tasks

$$470 \quad \theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(\phi_i, \mathcal{D}_i^{val}). \quad (1)$$

471 This enables MAML to find a good parameter initial-
 472 ization that can quickly converge to a minimum,
 473 making it suitable for Few-Shot Learning. The
 474 model was used for Few-Shot text classification
 475 (Han et al., 2018; Obamuyide and Vlachos, 2019;
 476 Jiang et al., 2018; Bao et al., 2020), where each
 477 class is considered a task. Additionally, Jiang
 478 et al. (2018) introduces task-agnostic parameters
 479 and task-specific parameters to MAML, which they
 480 call **ATAML**, outperforming vanilla MAML on
 481 Few-Shot topic classification. MAML has also
 482 seen applications in Few-Shot domain adaption,
 483 e.g. Few-Shot dialogue system (Lin et al., 2019;
 484 Mi et al., 2019; Qian and Yu, 2019), where each do-
 485 main dialog is treated as a task. One problem that
 486 MAML has, is that it is both computationally and
 487 memory intensive since it needs to calculate second
 488 derivatives in equation (1) as we get a nested back-
 489 propagation, where second derivatives may come
 490 up. **First-Order MAML (FOMAML)** and **REP-**
 491 **TILE** (Finn et al., 2017; Nichol et al., 2018) are
 492 methods which approximate the second derivative.
 493 One of the biggest challenges is to apply MAML
 494 to diverse tasks, as most applications are limited
 495 to similar train and test tasks, e.g. domain adap-
 496 tion tasks or to simulated classification datasets
 497 where each label is considered a task. Furthermore,
 498 even though the approach itself is model agnostic,
 499 meaning we can combine any model representation
 500 and any differentiable objective, the approach is
 501 restricted to tasks that have the same label space
 502 since to learn a good initialization, MAML requires
 503 sharing model parameters, including softmax clas-
 504 sification layers across tasks.

505 To enable MAML to learn across diverse tasks
 506 with disjoint label spaces, Bansal et al. (2019) pro-
 507 poses **LEOPARD**, which uses a parameter genera-
 508 tor, which learns on \mathcal{D}_i^{tr} to generate *task-dependent*
 509 *initial softmax classification parameters* for any
 510 specific task. Furthermore, the approach transforms
 511 the text input into a feature representation by us-
 512 ing a (shared) BERT model across tasks. To find
 513 a good parameter initialization, LEOPARD uses a

514 modified MAML-based adaptation method by dis-
 515 tinction between *task-specific parameters*, which
 516 are adapted per task, and *task-agnostic parameters*,
 517 which are shared across tasks. This is similar to
 518 Jiang et al. (2018). This allows for more efficient
 519 adaptation. Since BERT has a high number of pa-
 520 rameters, LEOPARD uses lower-layers of BERT as
 521 task-agnostic parameters and higher-level layer and
 522 the softmax generating function as task-specific pa-
 523 rameters. Since we already used \mathcal{D}_i^{tr} to generate
 524 task-dependent initial softmax classification param-
 525 eters, we use subsequent batches for adaption. On
 526 the contrary to vanilla MAML, LEOPARD can
 527 handle test tasks that are notably different from
 528 the training tasks. They evaluate LEOPARD using
 529 target tasks that were *not seen during training* and
 530 evaluate on their *entire test* set after finetuning on
 531 K examples per label from the corresponding train-
 532 ing set. The target tasks were selected such that
 533 they differ significantly from the training task and
 534 have a varying number of labels. They show that on
 535 average LEOPARD performs significantly better
 536 than the chosen baselines, BERT-base model (De-
 537 vlin et al., 2019), Multi-task BERT (comparable to
 538 Liu et al. (2019)) and a Prototypical Network (Snell
 539 et al., 2017) that uses BERT-base as the underlying
 540 neural model. With that experiment, they show that
 541 LEOPARD can leverage Meta-learning to learn
 542 a more general-purpose parameter initialization
 543 that can then be used to solve completely unseen
 544 new tasks with just a few examples. Furthermore,
 545 Bansal et al. (2019) evaluate Few-Shot Domain-
 546 transfer, showing that LEOPARD performs on par
 547 or better than the baselines. They also show that
 548 prototypical networks give competitive results on
 549 domain-transfer tasks. One disadvantage of LEOP-
 550 ARD is that it requires labeled data from many
 551 different tasks, for training and also hyperparam-
 552 eter tuning. Additionally, it suffers from overfitting
 553 to the training task-distribution (Bansal et al., 2019,
 554 2020) (*Meta-overfitting*), leaving room for a more
 555 efficient adaption to diverse tasks.

556 3.3 Discussion: Meta-Learning for NLP 557 Few-Shot Tasks

558 One of the main challenges in Meta-Learning (in
 559 general) is to create training tasks that enable Meta-
 560 learning algorithms to find a good initialization set
 561 to solve the target task (Vinyals et al., 2017). As
 562 previously mentioned, many applications create
 563 training tasks from a fixed task dataset, where we

564 have many labels, by subsampling from the set of
565 labels. While it enables to generalize to unseen
566 labels, this can also lead to overfitting to the train-
567 ing task distribution, making it hard to generalize
568 to unseen tasks (Yin et al., 2020a). Furthermore,
569 one of the reasons why most Meta-learning algo-
570 rithms were first proposed in image classification
571 problems is because they have big labeled sets with
572 a large number of labels. In NLP tasks, however,
573 they are often restricted to a small number of la-
574 bels, e.g. sentiment classification has only a few
575 discrete labels. To remedy this, Bansal et al. (2020)
576 propose a self-supervised approach to generate a
577 Meta-learning task distribution from an unlabeled
578 text by masking words from a specified vocabu-
579 lary (or subsets of it) and posing it as a multi-class
580 classification. Combining the generated tasks with
581 the available supervised tasks can improve Meta-
582 learning algorithms, such as LEOPARD (Bansal
583 et al., 2020). However, as these generated tasks
584 are only (masked language) classification tasks,
585 this can lead to a narrow training-task distribution.
586 Additionally, most of the research only explores
587 classification problems, leaving room for future
588 work to expand into more diverse problem struc-
589 tures and to find more suitable ways to generate
590 diverse Meta-learning tasks.

591 One major obstacle for Meta-learning ap-
592 proaches is to solve diverse NLP Few-Shot tasks.
593 Meta-Learning approaches may work well for sim-
594 ulated datasets, where we just subsample labels
595 from one single task dataset and define them as
596 training tasks because the underlying task does not
597 change in this situation, e.g. the model was “pre-
598 trained” to solve translation tasks. However, if you
599 want to test on a truly unseen task, the model has
600 to first learn the underlying task from a few given
601 examples. Jiang et al. (2018); Bansal et al. (2019)
602 mitigate this problem by introducing task-specific
603 parameters and task-agnostic parameters for more
604 efficient adaption. Another interesting approach
605 for future work could be to combine Meta learning
606 with additional task information, e.g. task descrip-
607 tions, to solve new diverse tasks (approaches in
608 Section 2 do this).

609 4 K -Shot Cross-Lingual Transfer with 610 Multilingual Language Models

611 This section will deal with K -Shot Cross-Lingual
612 Transfer as a use-case of Few-Shot Learning.
613 Achieving SOTA on (monolingual) NLP tasks is

614 usually done by using transformers, pre-trained on
615 language modeling objectives in a semi-supervised
616 fashion, and then finetuning on a specific NLP task,
617 which in return need a lot of labeled training data.
618 Those are available in common languages, such
619 as the English language, however, in low resource
620 languages models fail to generalize well. The idea
621 is to transfer the knowledge about a task from a
622 high resource language to another low resource
623 language, called cross-lingual transfer (CLT).

624 To achieve CLT between tasks from different
625 languages, one has to induce a shared repre-
626 sentation space between the source and target
627 language. Previous SOTA methods used to
628 induce continuous cross-lingual representation
629 spaces by using cross-lingual word embeddings
630 (Mikolov et al., 2013b; Glavaš et al., 2019) and
631 sentence embeddings (Artetxe and Schwenk,
632 2019). However, with transformers getting popular,
633 this survey will focus on inducing multilingual
634 word embeddings with transformers.
635

636 4.1 Zero-Shot Cross-Lingual Transfer

637 One way to try to combat sparsely labeled train-
638 ing data in one language is by pretraining trans-
639 former models on multiple languages and auto-
640 matically induce a multilingual word embedding.
641 This idea gave rise to powerful massively multi-
642 lingual transformers, such as mBert, XLM-R, and
643 the recently introduced mT5 (Devlin et al., 2019;
644 Conneau et al., 2020; Xue et al., 2021). These archi-
645 tectures can encode text from any of the languages
646 seen in pretraining and allows for a very straightfor-
647 ward approach to Zero-Shot cross-lingual model
648 transfer: Finetune the model using task-specific
649 supervised training data from one high resource
650 language (*source-training*) and predict on other
651 languages by feeding the target language text into
652 the finetuned model. Pires et al. (2019) show ef-
653 fective results of Zero-Shot cross-lingual transfer
654 with mBERT on POS tagging and NER for related
655 languages. Furthermore, Wu et al. (2020); K et al.
656 (2020) show the cross-lingual potential of mBERT
657 by extending the analysis. Nevertheless, the litera-
658 ture mostly showed good results in languages that
659 were from the same language family or that had
660 a large corpus in pretraining, languages such as
661 German, Spanish or French. This concern is raised
662 by multiple sources (Lauscher et al., 2020; Wu and
663 Dredze, 2020), which show that the performance

664 drops huge for distant target languages and tar- 714
665 get languages that have small pre-training corpus. 715
666 Furthermore, Lauscher et al. (2020) empirically 716
667 show that for massively multilingual transform- 717
668 ers, pre-training corpora sizes affect the Zero-Shot 718
669 performance in higher-level language understand- 719
670 ing tasks (e.g. NLI and QA), whereas the results 720
671 in lower-level language understanding tasks are 721
672 more impacted by typological language proxim- 722
673 ity. To summarize, Zero-Shot cross-lingual transfer 723
674 with source training is effective for languages that 724
675 are linguistically similar and languages that have a 725
676 great amount of data for pre-training. However, this 726
677 scenario is almost always never the case for low 727
678 resource languages, where cross-lingual transfer is 728
679 needed. The next section will investigate Few-Shot 729
680 transfer to mitigate the transfer gap. 730

681 4.2 Few-Shot Cross-Lingual Transfer

682 To improve upon the results of Zero-Shot CLT, 732
683 which only uses source training, we now addi- 733
684 tionally exploit the K task-specific examples in 734
685 the target language (Few-Shot cross-lingual sce- 735
686 nario) by further finetuning on those K examples 736
687 (*target-adapting*). Lauscher et al. (2020) experi- 737
688 ment with Few-Shot CLT on lower-level structured 738
689 prediction tasks (POS tagging, dependency parsing, 739
690 and NER) and higher-level language understanding 740
691 tasks (NLI and QA) with varying numbers of K 741
692 examples. They show that distant languages gain 742
693 much more in performance from Few-Shot data 743
694 than closely related languages. Hedderich et al. 744
695 (2020) use Few-Shot CLT on NER task on genuine 745
696 low-resource languages like Hausa and isiXhosa, 746
697 also showing significant improvements by finetun- 747
698 ing on the few examples. Zhao et al. (2020) applied 748
699 Few-Shot CLT with mBERT on POS, NER, and 749
700 sequence classification, observing the same phe- 750
701 nomenon. In summary, additional finetuning on 751
702 the given few examples from the target language 752
703 can significantly improve performances on distant 753
704 languages - Exactly where Zero-Shot CLT fails. 754
705 Since we only have to finetune on a small set of 755
706 examples, this additional finetuning is not compu- 756
707 tationally expensive but shows promising results. 757

708 As we only discussed "naively" finetuning for 758
709 target adaption, one could further investigate how 759
710 to exploit the given examples efficiently. Zhao et al. 760
711 (2020) investigated freezing parameters during fine- 761
712 tuning to mitigate the overfitting problem, however, 762
713 experiments show no significant improvements in 763

performance. To use the few given examples more 714
efficiently, Nooralahzadeh et al. (2020) use MAML 715
to further find optimal initialization parameters (af- 716
ter source training), which then can be used for 717
either Zero-Shot or again finetuning in a Few-Shot 718
setup. However, the method requires many train- 719
ing tasks in low-resource languages. Future work 720
could focus on using Meta-Learning further. 721

722 One downside of all Few-Shot CLT approaches 722
is that you need labeled data in the low resource 723
target language, which is typically hard to acquire. 724
It may become costly to annotate data for minor 725
languages, however as Lauscher et al. (2020) show, 726
even 10 annotated instances can give substantial 727
performance improvement. This begs the question 728
if annotating data is more cost-efficient in the long 729
run than using GPU hours. 730

731 5 Conclusion and Discussion

732 We studied two methods to tackle Few-Shot tasks 732
in NLP: Using pre-trained Language Models and 733
Meta-learning. Even though Meta-Learning pro- 734
vides diverse applications as most methods are task 735
and model agnostic, they struggle to solve unseen 736
diverse NLP tasks. Future work should investi- 737
gate how to improve generalization to new tasks. 738
Pre-trained language models can be effective by 739
reformulating NLP tasks as language model prob- 740
lems, enabling Few-Shot abilities. However these 741
methods require manual work to find a good re- 742
formulation and they favor tasks, that can be nat- 743
urally reformulated as a "fill-in-the-blank" task. 744
We then discussed a use-case of Few-Shot Learn- 745
ing: Few-Shot CLT. In CLT, we have the chance 746
to first finetune in a rich-resource language, and 747
then transfer the knowledge to a low-resource lan- 748
guage. Using more sophisticated methods to train 749
on high resource languages, e.g. Meta-Learning 750
(Nooralahzadeh et al., 2020), can improve perfor- 751
mance and is a promising research direction. Nev- 752
ertheless, most methods need labeled examples in 753
low resource languages, making them expensive 754
to obtain. As previously discussed in Section 1.3, 755
almost all Few-Shot techniques have high variance. 756
Therefore, we identify the necessity of standardiza- 757
tion of Few-Shot datasets. As a final word, there 758
are other approaches to Few-Shot Learning in NLP 759
that was not discussed in this survey, e.g. unify- 760
ing NLP tasks formats (McCann et al., 2018b; Yin 761
et al., 2020b; Raffel et al., 2020; Khashabi et al., 762
2020). 763

764
765
766
767

768
769
770
771

772
773
774
775

776
777
778

779
780
781

782
783

784
785
786
787
788
789
790

791
792
793
794
795
796
797
798
799
800
801
802

803
804
805

806
807
808
809
810

811
812

813
814
815
816
817

References

Mikel Artetxe and Holger Schwenk. 2019. [Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond](#).

Trapit Bansal, Rishikesh Jha, and Andrew McCallum. 2019. [Learning to few-shot learn across diverse natural language classification tasks](#). *CoRR*, abs/1911.03863.

Trapit Bansal, Rishikesh Jha, Tsendsuren Munkhdalai, and A. McCallum. 2020. [Self-supervised meta-learning for few-shot natural language classification tasks](#). *ArXiv*, abs/2009.08445.

Yujia Bao, Menghua Wu, Shiyu Chang, and Regina Barzilay. 2020. [Few-shot text classification with distributional signatures](#).

Yonatan Belinkov and Yonatan Bisk. 2018. [Synthetic and natural noise both break neural machine translation](#).

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. [Longformer: The long-document transformer](#).

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does bert look at? an analysis of bert’s attention](#).

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#).

Andrew M. Dai and Quoc V. Le. 2015. [Semi-supervised sequence learning](#).

Shumin Deng, Ningyu Zhang, Jiaojian Kang, Yichi Zhang, Wei Zhang, and Huajun Chen. 2020. [Meta-learning with dynamic-memory-based prototypical network for few-shot event detection](#). In *Proceedings of the 13th International Conference on*

Web Search and Data Mining, WSDM ’20, page 151–159, New York, NY, USA. Association for Computing Machinery. 818
819
820

J. Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*. 821
822
823
824

Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. 2020. [Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping](#). 825
826
827
828
829

Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. [Model-agnostic meta-learning for fast adaptation of deep networks](#). *CoRR*, abs/1703.03400. 830
831
832

Karèn Fort. 2016. *Collaborative Annotation for Reliable Natural Language Processing: Technical and Sociological Aspects*, 1st edition. Wiley-IEEE Press. 833
834
835
836

Tianyu Gao, Adam Fisch, and Danqi Chen. 2020. [Making pre-trained language models better few-shot learners](#). 837
838
839

Tianyu Gao, Xu Han, Zhiyuan Liu, and Maosong Sun. 2019. [Hybrid attention-based prototypical networks for noisy few-shot relation classification](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):6407–6414. 840
841
842
843
844

Goran Glavaš, Robert Litschko, Sebastian Ruder, and Ivan Vulić. 2019. [How to \(properly\) evaluate cross-lingual word embeddings: On strong baselines, comparative analyses, and some misconceptions](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 710–721, Florence, Italy. Association for Computational Linguistics. 845
846
847
848
849
850
851
852

Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. 2021. [Supervised contrastive learning for pre-trained language model fine-tuning](#). 853
854
855

Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018. [FewRel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4803–4809, Brussels, Belgium. Association for Computational Linguistics. 856
857
858
859
860
861
862
863

Michael A. Hedderich, David Adelani, Dawei Zhu, Jesujoba Alabi, Udia Markus, and Dietrich Klakow. 2020. [Transfer learning and distant supervision for multilingual transformer models: A study on african languages](#). 864
865
866
867
868

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. [Distilling the knowledge in a neural network](#). 869
870

871	Bei Hui, Liang Liu, J. Chen, X. Zhou, and Yuhui Nian. 2020. Few-shot relation classification by context attention-based prototypical networks with bert. <i>EURASIP Journal on Wireless Communications and Networking</i> , 2020:1–17.	Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms .	923 924
876	Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems .	Farhad Nooralahzadeh, Giannis Bekoulis, Johannes Bjerva, and Isabelle Augenstein. 2020. Zero-shot cross-lingual transfer with meta learning .	925 926 927
878	Xiang Jiang, Mohammad Havaei, G. Chartrand, H. Chouaib, Thomas Vincent, Andrew Jesson, Nicolas Chapados, and S. Matwin. 2018. Attentive task-agnostic meta-learning for few-shot text classification.	Abiola Obamuyide and Andreas Vlachos. 2019. Model-agnostic meta-learning for relation classification with limited supervision . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5873–5879, Florence, Italy. Association for Computational Linguistics.	928 929 930 931 932 933
883	Karthikeyan K, Zihan Wang, Stephen Mayhew, and Dan Roth. 2020. Cross-lingual ability of multilingual bert: An empirical study .	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations .	934 935 936 937
886	Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hananeh Hajishirzi. 2020. UNIFIEDQA: Crossing format boundaries with a single QA system . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1896–1907, Online. Association for Computational Linguistics.	Jason Phang, Thibault Févry, and Samuel R. Bowman. 2019. Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks .	938 939 940
893	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations .	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual bert?	941 942
897	Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. 2020. From zero to hero: On the limitations of zero-shot cross-lingual transfer with multilingual transformers .	Kun Qian and Zhou Yu. 2019. Domain adaptive dialog generation via meta learning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2639–2649, Florence, Italy. Association for Computational Linguistics.	943 944 945 946 947
901	Zhaojiang Lin, Andrea Madotto, Chien-Sheng Wu, and Pascale Fung. 2019. Personalizing dialogue agents via meta-learning .	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	948 949 950 951
904	Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019. Multi-task deep neural networks for natural language understanding .	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer .	952 953 954 955 956
907	Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. 2018a. Learned in translation: Contextualized word vectors .	Sebastian Ruder. 2019. <i>Neural Transfer Learning for Natural Language Processing</i> . Ph.D. thesis, National University of Ireland, Galway.	957 958 959
910	Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018b. The natural language decathlon: Multitask learning as question answering .	Timo Schick and Hinrich Schütze. 2021a. Exploiting cloze questions for few shot text classification and natural language inference .	960 961 962
914	Fei Mi, Minlie Huang, Jiyong Zhang, and Boi Faltings. 2019. Meta-learning for low-resource natural language generation in task-oriented dialogue systems .	Timo Schick and Hinrich Schütze. 2021b. It’s not just size that matters: Small language models are also few-shot learners .	963 964 965
917	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space .	J. Snell, Kevin Swersky, and R. Zemel. 2017. Prototypical networks for few-shot learning. In <i>NIPS</i> .	966 967
920	Tomas Mikolov, Quoc V. Le, and Ilya Sutskever. 2013b. Exploiting similarities among languages for machine translation .	Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip H. S. Torr, and Timothy M. Hospedales. 2018. Learning to compare: Relation network for few-shot learning .	968 969 970 971
922		Derek Tam, R. R. Menon, M. Bansal, Shashank Srivastava, and Colin Raffel. 2021. Improving and simplifying pattern exploiting training . <i>ArXiv</i> , abs/2103.11955.	972 973 974 975

976	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	Mengjie Zhao, Yi Zhu, Ehsan Shareghi, Roi Reichart,	1029
977	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	Anna Korhonen, and Hinrich Schütze. 2020. A	1030
978	Kaiser, and Illia Polosukhin. 2017. Attention is all	closer look at few-shot crosslingual transfer: Vari-	1031
979	you need.	ance, benchmarks and baselines.	1032
980	Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Ko-		
981	ray Kavukcuoglu, and Daan Wierstra. 2017. Match-		
982	ing networks for one shot learning.		
983	Alex Wang, Yada Pruksachatkun, Nikita Nangia,		
984	Amanpreet Singh, Julian Michael, Felix Hill, Omer		
985	Levy, and Samuel R. Bowman. 2020. Superglue: A		
986	stickier benchmark for general-purpose language un-		
987	derstanding systems.		
988	Alex Wang, Amanpreet Singh, Julian Michael, Felix		
989	Hill, Omer Levy, and Samuel R. Bowman. 2019.		
990	Glue: A multi-task benchmark and analysis platform		
991	for natural language understanding.		
992	Shijie Wu, Alexis Conneau, Haoran Li, Luke Zettle-		
993	moyer, and Veselin Stoyanov. 2020. Emerging cross-		
994	lingual structure in pretrained language models.		
995	Shijie Wu and Mark Dredze. 2020. Are all lan-		
996	guages created equal in multilingual bert? In		
997	RepLANLP@ACL.		
998	Linting Xue, Noah Constant, Adam Roberts, Mi-		
999	hir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya		
1000	Barua, and Colin Raffel. 2021. mt5: A massively		
1001	multilingual pre-trained text-to-text transformer.		
1002	Mingzhang Yin, George Tucker, Mingyuan Zhou,		
1003	Sergey Levine, and Chelsea Finn. 2020a. Meta-		
1004	learning without memorization.		
1005	Wenpeng Yin, Nazneen Fatema Rajani, Dragomir		
1006	Radev, Richard Socher, and Caiming Xiong. 2020b.		
1007	Universal natural language processing with limited		
1008	annotations: Try few-shot textual entailment as a		
1009	start.		
1010	Dani Yogatama, Cyprien de Masson d’Autume, Jerome		
1011	Connor, Tomas Kocisky, Mike Chrzanowski, Ling-		
1012	peng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu,		
1013	Chris Dyer, and Phil Blunsom. 2019. Learning and		
1014	evaluating general linguistic intelligence.		
1015	Mo Yu, Xiaoxiao Guo, Jinfeng Yi, Shiyu Chang, Saloni		
1016	Potdar, Yu Cheng, Gerald Tesauro, Haoyu Wang,		
1017	and Bowen Zhou. 2018. Diverse few-shot text clas-		
1018	sification with multiple metrics. In <i>Proceedings of</i>		
1019	<i>the 2018 Conference of the North American Chap-</i>		
1020	<i>ter of the Association for Computational Linguistics:</i>		
1021	<i>Human Language Technologies, Volume 1 (Long Pa-</i>		
1022	<i>pers)</i> , pages 1206–1215, New Orleans, Louisiana.		
1023	Association for Computational Linguistics.		
1024	Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q.		
1025	Weinberger, and Yoav Artzi. 2021. Revisiting few-		
1026	sample bert fine-tuning.		
1027	Yu Zhang and Qiang Yang. 2021. A survey on multi-		
1028	task learning.		