

A Survey on Stability of Learning with Limited Labelled Data and its Sensitivity to the Effects of Randomness

BRANISLAV PECHER*, Faculty of Information Technology, Brno University of Technology, Czechia

IVAN SRBA, Kempelen Institute of Intelligent Technologies, Slovakia

MARIA BIELIKOVA†, Kempelen Institute of Intelligent Technologies, Slovakia

Learning with limited labelled data, such as prompting, in-context learning, fine-tuning, meta-learning or few-shot learning, aims to effectively train a model using only a small amount of labelled samples. However, these approaches have been observed to be excessively sensitive to the effects of uncontrolled randomness caused by non-determinism in the training process. The randomness negatively affects the stability of the models, leading to large variances in results across training runs. When such sensitivity is disregarded, it can unintentionally, but unfortunately also intentionally, create an imaginary perception of research progress. Recently, this area started to attract research attention and the number of relevant studies is continuously growing. In this survey, we provide a comprehensive overview of 415 papers addressing the effects of randomness on the stability of learning with limited labelled data. We distinguish between four main tasks addressed in the papers (investigate/evaluate; determine; mitigate; benchmark/compare/report randomness effects), providing findings for each one. Furthermore, we identify and discuss seven challenges and open problems together with possible directions to facilitate further research. The ultimate goal of this survey is to emphasise the importance of this growing research area, which so far has not received an appropriate level of attention, and reveal impactful directions for future research.

CCS Concepts: • **General and reference** → **Surveys and overviews**; • **Computing methodologies** → *Artificial intelligence; Machine learning; Learning paradigms.*

Additional Key Words and Phrases: randomness, stability, sensitivity, meta-learning, large language models, fine-tuning, prompting, in-context learning, instruction-tuning, prompt-based learning, PEFT, literature survey

ACM Reference Format:

Branislav Pecher, Ivan Srba, and Maria Bielikova. 2024. A Survey on Stability of Learning with Limited Labelled Data and its Sensitivity to the Effects of Randomness. *ACM Comput. Surv.* <https://doi.org/10.1145/3691339>

1 Introduction

The approaches for *learning with limited labelled data* are designed to achieve high performance in machine learning models even with few labels available [76, 134]. Under the term *learning with limited labelled data*, we understand any approach that is designed to work with a lack of labels, without any constraint on how the labelled samples are distributed, i.e., whether all the samples are from a single task or are distributed across different tasks. Although similar to the

* Also with Kempelen Institute of Intelligent Technologies.

† Also with slovak.AI.

Authors' Contact Information: [Branislav Pecher](mailto:branislav.pecher@kinit.sk), Faculty of Information Technology, Brno University of Technology, Brno, Czechia, branislav.pecher@kinit.sk; [Ivan Srba](mailto:ivan.srba@kinit.sk), Kempelen Institute of Intelligent Technologies, Bratislava, Slovakia, ivan.srba@kinit.sk; [Maria Bielikova](mailto:maria.bielikova@kinit.sk), Kempelen Institute of Intelligent Technologies, Bratislava, Slovakia, maria.bielikova@kinit.sk.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 1557-7341/2024/1-ART1

<https://doi.org/10.1145/3691339>

notion of few-shot learning, it represents a broader scope encompassing a larger number of possible approaches (some of which are incorrectly categorised as few-shot learning in current practice).

To deal with the limited labels, these approaches utilise additional information from different sources [19, 134], such as transferring knowledge from similar tasks and datasets. In the NLP domain, *prompting* and *in-context learning* (also called few-shot prompting) have recently emerged. In these techniques, a large pre-trained language model is “prompted” to predict a label for a single test sample by presenting it with task instructions and the test sample (and concatenation of a few labelled samples when *in-context learning* is used), without requiring any parameter update [76]. In addition, it is common to use *fine-tuning*, where the parameters, or their subset using parameter-efficient fine-tuning (PEFT) methods, of the pre-trained large language model are updated to optimise the model for the specific downstream task using only a few labelled samples [25, 33, 94]. Finally, *meta-learning* can be used, where the model is explicitly trained to quickly adapt to a new task with only a handful of examples by learning how best to learn across a large number of related tasks with few labelled samples each [2, 49].

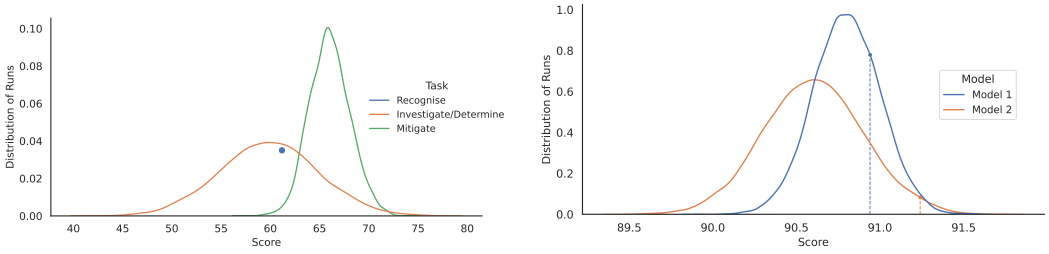
However, a significant problem observed for these approaches is their sensitivity to the effects of uncontrolled randomness, which negatively affects their stability. Under *stability* and its opposite term *sensitivity*, we understand a property of a learning algorithm or a model that indicates what influence the *small-scale and random changes* (or perturbations) in the input *data* and *parameters* have on its outputs. Such random changes are introduced by different *randomness factors* that represent the non-deterministic decisions in the training process, such as random initialisation of model parameters or data shuffling [44, 110]. These *randomness factors* represent the main point around which the effects of randomness are addressed.

The uncontrolled effects of randomness can have a massive impact on the stability of the utilised machine learning approaches. Running the training multiple times on the same dataset, with the same setup and hyperparameters, may lead to large deviations in the final performance [2, 33, 83, 94]. Changing only the order of samples or answers in multi-choice question answering when using in-context learning can lead the model from state-of-the-art predictions to random guesses [83, 149, 191]. Choosing a different set of adaptation data in meta-learning can lead to a difference between minimum and maximum performance being up to 90% [2]. Even though the effects of randomness are present in all cases of machine learning, their impact on stability is especially significant in low-data regimes.

If the uncontrolled randomness is not appropriately addressed in such low-data regimes, it can have significant, non-negligible negative consequences. In comparisons and benchmarks, changing only the random seed may lead to completely different model rankings [7, 86]. It may prohibit an objective performance comparison of newly designed methods to the state-of-the-art baselines, as it makes it difficult to conclude if a specific change made a meaningful difference or if the model just “got lucky” (especially when the difference between their performance is similar to the deviation caused by the randomness). In such cases, a method can be incorrectly denoted as state-of-the-art only based on a more favourable random chance [120] (as illustrated in Figure 1b). The uncontrolled randomness can unintentionally, but unfortunately also intentionally (by cherry-picking), create an imaginary perception of research progress. As a result, researchers even suggest that we should be very vigilant when relying on models and approaches that are significantly affected by sensitivity (e.g., large language models) [191]. Finally, randomness has been identified as a significant obstacle that negatively affects reproducibility [5, 23].

The effects of randomness are addressed in various depths - from a pure *recognition* of the effects, through a deeper *investigation* and *determining* the origin of the effects, up to their partial or full *mitigation* (illustrated in the Figure 1). First, the effects are simply recognised, resulting in the reporting of a more representative, albeit single, aggregated value from multiple training and

evaluation runs (denoted as “Recognise”). As no further analysis is provided, this still leads to a biased comparison between approaches, as one approach may show better results only due to random chance and unintentional cherry-picking (illustrated in Figure 1b). Another group of works perform analyses of the effects of randomness, estimate the distribution of results from multiple runs and compare the distributions (denoted as “Investigate/Determine”). The effects of the randomness are often analysed in more detail in order to determine the real origin of randomness in the training process (e.g., under-specification, which causes slightly different parameter initialisation to converge to different minima). As the distribution is only estimated, the deviation is still not reduced. The effects of randomness are fully addressed when they are mitigated, reducing the deviation in the distribution of results (denoted as “Mitigate”). However, effective mitigation requires an understanding of the effects, their importance, and the origin of randomness provided by the analysis of the effects of randomness (provided by the “Investigate/Determine” task).



(a) Various depths the effects of randomness can be addressed in the papers

(b) Comparison based on a single value and based on distribution of the results, inspired by [120]

Fig. 1. (a) The effects of randomness can be addressed in various depths in the papers. (b) If not taken into consideration, randomness can introduce bias into comparison results, causing one approach to show better results only due to random chance and unintentional cherry-picking.

In this paper, we conduct a comprehensive survey of 415 papers that address the effects of randomness, either by simply recognising the effects of randomness on stability ($N=254$), or by investigating/determining the characteristics of these effects in more detail and/or exploring the appropriate means of their mitigation ($N=161$)¹. We provide a more in-depth analysis of the 161 papers that address the effects of randomness in more detail. Finally, we aggregate findings from this analysis and use them to identify seven challenges and open problems related to the study of the effects of randomness, while providing future directions on how to address them. For the purpose of this survey, we follow the PRISMA methodology [93, 100] to systematically identify the papers in the defined scope.

This is the first survey of its kind that provides a comprehensive and systematic literature review specifically focused on the impact of the effects of randomness on stability across various approaches for dealing with limited labelled data. As opposed to previous papers that provide an overview of possible *randomness factors* and investigate these effects [44, 110, 134], provide an overview for a single task (e.g., mitigating the effects) [158], or only recognise the problem randomness as an open problem [76, 134], we focus on all tasks for addressing the effects of randomness, including investigation, but also determining their origin, mitigation, and benchmarking in the presence of these effects of randomness.

The purpose of this survey is to emphasise the importance of this research area, which so far has not received the appropriate level of attention from researchers and practitioners. First, it should

¹The digital appendix with a full list of papers is available at <https://kinit.sk/public/acm-csur-sensitivity-survey.html>

serve existing or new researchers who are explicitly tackling the topic of randomness and its effects on the stability of *learning with limited labelled data* to support their research (by providing an overview of the state-of-the-art, identification of open problems and future directions). Secondly, its purpose is to inform researchers and practitioners utilising the *learning with limited labelled data* about the consequences of unaddressed randomness and how to effectively prevent them. To better achieve these purposes, we summarised all the identified papers along with their categorisation in the [digital appendix of this survey](#) (which is further described in Section A of the supplementary material).

The rest of the paper is structured as follows. In Section 2, we describe the scope of the survey, the difference to existing surveys and the methodology used for identifying relevant papers. Taxonomy applied to analyse and categorise the identified papers is described in Section 3. Afterwards, we describe the different tasks for addressing the effects of randomness, categorised based on our taxonomy, in Sections 4, 5, 6 and 7, with a summary of the findings at the end of each such section. The identified challenges and open problems are described in Section 8. Finally, the survey is concluded in Section 9 with a summary of our contributions and findings.

2 Background

2.1 Scope of the Survey

In this survey, we focus on a specific part of *stability* and *sensitivity* of machine learning approaches. We focus on papers that explore the effects of small perturbations to **multiple inputs** of the training process on the final performance after the training process is done. A similar, but at the same time quite distinct property is the *robustness* of the model. It deals with a more large scale or systematic changes in the input, such as distribution shift or the presence of adversarial examples. As the *adversarial robustness* is already extensively studied in the context of learning with limited labelled and the effects are often systematic, we consider it to be out of scope for this survey. Finally, we view the *stability* to be a subset of the *reproducibility* of the model. Besides randomness in the training process, reproducibility is impacted by other factors such as selective reporting, biases, data and source code availability, and many others as specified in [44]. These additional factors are out of scope for this survey.

Out of a wide spectrum of machine learning approaches, we specifically focus on *learning with limited labelled data*, which is designed to train or prompt a model with only a limited number of labelled samples (typically 0-50 per class). Not requiring large sets of labelled data, these approaches are usable across a broader set of domains, making them also more popular. Due to this popularity, they also cover almost the whole research area of addressing the effects of randomness on stability, as there are almost no other works focusing on approaches that utilise information from unlabelled sets of data. Our focus is on five prevalent groups of approaches: 1) *meta-learning*; 2) *language model (LM) fine-tuning*; 3) *prompting/in-context learning*; 4) *prompt-based learning*; and 5) *parameter-efficient fine-tuning (FT)*. Even though we do not explicitly focus on a specific modality, the popularity of the approaches impacts the scope of this survey. The most prominent modality is text, followed by images, as the majority of the papers focus on language model *fine-tuning* and *in-context learning*. Finally, even though we focus on the setting of limited data, we also include works that are relevant for this survey that utilise large sets of data. However, due to the search query used, such papers are not systematically covered and may be incomplete (even though we specifically designed one step in the methodology to catch these papers).

2.2 Existing Surveys

To the best of our knowledge, there are no existing surveys addressing the effects of randomness on stability across all tasks ranging from recognise to mitigate. However, there are relevant surveys on related topics that partially touch on the topics of *stability* or *sensitivity*. First of all, some papers provide an overview of possible sources of randomness, although with a focus on replicability [44]. In some surveys on *learning with limited labelled data* and/or *few-shot learning*, the sensitivity to the effects of randomness is already mentioned as a problem that should be addressed [34, 57, 76, 134], but no additional details are provided. Finally, the most recent relevant survey [158] focuses on sample selection strategies for mitigating in-context learning sensitivity to sample choice. While relevant, it provides a significantly narrower perspective than the one given by this comprehensive survey.

To fill this gap, we aim to conduct the first comprehensive overview of the effects of randomness on the stability of learning with limited labelled data and provide novel perspectives that are not easily visible from the standard analysis of related works.

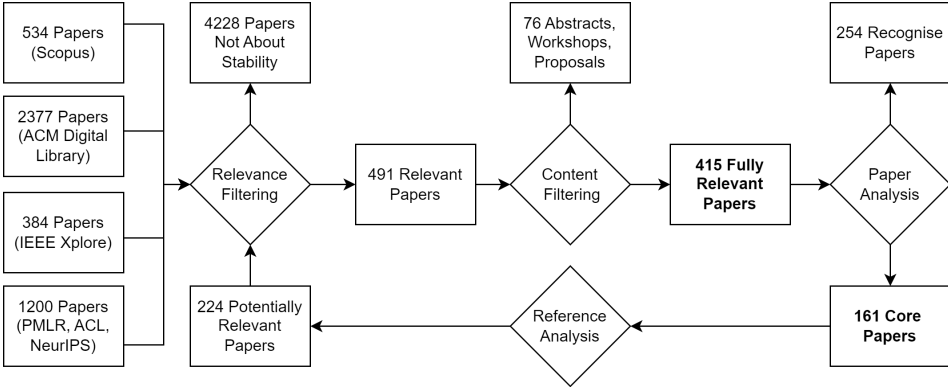


Fig. 2. Process for identifying and categorising papers for this survey. Due to strong clustering effects, additional papers are identified using a reference analysis on the most relevant identified papers, i.e., including papers cited in them, as well as the ones that cite them.

2.3 Survey Methodology

To systematically identify the set of relevant papers for our comprehensive survey, we follow the PRISMA methodology [93, 100], performing the following steps (as illustrated in Figure 2):

- (1) **Identification of relevant papers.** A keyword search query is defined based on the related terms for our defined scope identified in different papers. Multiple digital libraries are searched using the defined query to identify a set of potentially relevant papers.
- (2) **Relevance filtering.** The potentially relevant papers are categorised into irrelevant and relevant based on the context the defined keywords are used in.
- (3) **Content filtering.** The non-full papers or papers released in unrelated venues, such as extended abstracts, proposals for talks or papers appearing in unrelated workshops not dealing with stability, are filtered out.
- (4) **Paper analysis.** The papers are analysed based on their depth of addressing the randomness effects. The papers that only recognise the problem of stability without further investigation are referred to as *recognise* papers. The remaining papers are referred to as *core* papers.

- (5) **Reference analysis.** To discover additional potentially relevant papers (especially those relevant for the survey that do not necessarily deal with limited data), the papers cited in the *core* papers; as well as papers citing the *core* papers, are examined and their relevance is determined based on their title and abstract. The papers deemed potentially relevant are an additional input into the second step of this process.
- (6) **Analysis and categorisation.** The identified *core* papers are analysed and categorised based on the taxonomy (as defined in Section 3) to provide a comprehensive overview and to identify challenges and open problems.

As a result, 415 papers are identified — out of them, 254 papers are classified as *recognise* papers, and 161 as *core* papers. The number of identified papers grouped by the publication year is presented in Figure 3. The categorisation of the selected *core* papers is showcased in Table 1. Moreover, the full list of all identified papers (including *recognise* papers), with their categorisation and additional metadata is available in the [digital appendix](#). For a more detailed implementation of the survey methodology, see the supplementary material (Section C).

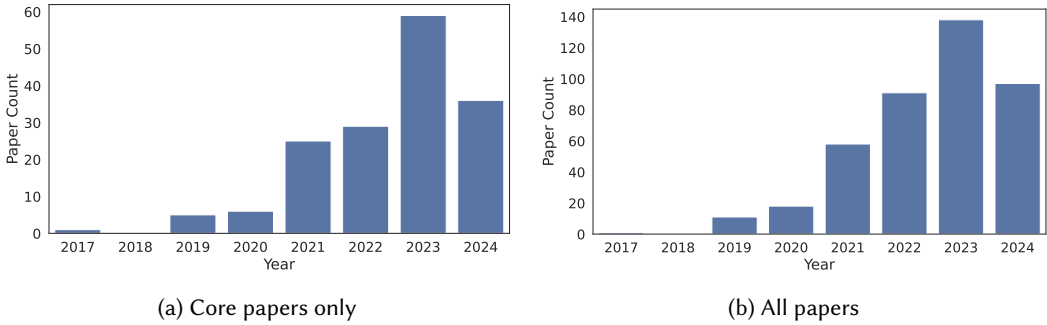


Fig. 3. Number of papers dealing with randomness grouped by the year. Figure 3a shows only the *core* papers that focus on addressing the effects of randomness in more detail, while Figure 3b also includes the papers that only *recognise* the problem. The problem started to attract attention only in recent years. Year 2024 covers papers published until July 2024.

3 Taxonomy for Literature Analysis and Categorisation

The identified *core* papers are categorised based on multiple dimensions that describe their characteristics of interest. We define the following three separate primary properties that are used to categorise the papers: 1) what *tasks are performed to address the randomness*; 2) *randomness factors* addressed by the papers; and 3) what *machine learning (ML) approach* the papers use. The *performed tasks* property is the main property based on which we organise the rest of the survey, while the *randomness factors* and *machine learning approach* properties are mostly used to define the scope in this survey (see Section 2.1). We provide the paper distribution across individual *randomness factors* and specific *machine learning approaches* for all *core* papers in Table 2 (the paper distributions for specific tasks are included in the [digital appendix](#)).

3.1 Tasks Performed to Address the Randomness

This property specifies the different tasks for addressing the effects of randomness that are performed in the papers (the relations between these tasks, along with their input and output are illustrated in Figure 4). For this property, we define the following 7 values:

Table 1. A showcase of the selected *core* papers and their categorisation based on the proposed taxonomy. The full categorisation, along with the *recognise* papers is available in the [digital appendix](#).

References	ML Approach	Randomness Factor	Performed Tasks	Modality	Out-Of-Distribution	
					Interactions	x
[16]	✓		✓	✓	✓	x
[2]	✓		✓	✓	✓	x
[15]	✓		✓	✓	✓	x
[83]		✓	✓	✓	✓	x
[177]	✓		✓	✓	✓	x
[99]	✓		✓	✓	✓	x
[91]		✓	✓	✓	✓	x
[20]	✓		✓	✓	✓	x
[126]	✓	✓	✓	✓	✓	x
[56]	✓	✓	✓	✓	✓	x
[152]	✓	✓	✓	✓	✓	x
[160]	✓	✓	✓	✓	✓	x
[180]	✓	✓	✓	✓	✓	x
[130]	✓	✓	✓	✓	✓	x
[35]	✓	✓	✓	✓	✓	x
[33]	✓	✓	✓	✓	✓	x
[129]	✓	✓	✓	✓	✓	x
[94]	✓	✓	✓	✓	✓	x
[61]	✓	✓	✓	✓	✓	x
[120]	✓	✓	✓	✓	✓	x
[173]	✓	✓	✓	✓	✓	x
[75]	✓	✓	✓	✓	✓	x
[128]	✓	✓	✓	✓	✓	x
[89]	✓	✓	✓	✓	✓	x
[187]	✓	✓	✓	✓	✓	x
[17]	✓	✓	✓	✓	✓	x
[32]	✓	✓	✓	✓	✓	x
[184]	✓	✓	✓	✓	✓	x
[31]	✓	✓	✓	✓	✓	x
[30]	✓	✓	✓	✓	✓	x
[59]	✓	✓	✓	✓	✓	x
[53]	✓	✓	✓	✓	✓	x
[48]	✓	✓	✓	✓	✓	x
[110]	✓	✓	✓	✓	✓	x
[28]	✓	✓	✓	✓	✓	x
[111]	✓	✓	✓	✓	✓	x
[145]	✓	✓	✓	✓	✓	x
[169]	✓	✓	✓	✓	✓	x
[156]	✓	✓	✓	✓	✓	x
[90]	✓	✓	✓	✓	✓	x
[174]	✓	✓	✓	✓	✓	x
[182]	✓	✓	✓	✓	✓	x
[164]	✓	✓	✓	✓	✓	x
[18]	✓	✓	✓	✓	✓	x
[96]	✓	✓	✓	✓	✓	x
[36]	✓	✓	✓	✓	✓	x

(1) *Investigate* - the effects of randomness from different *randomness factors* are investigated using specifically designed experiments. For example, the training and evaluation is done multiple times to determine the distribution of performance across multiple runs.

Table 2. Relation between the *randomness factors* and *machine learning approaches* across all analysed papers. We can observe different popularity of the *randomness factors* across different *machine learning approaches*. The majority focus is denoted in **bold**. For example, *fine-tuning* mostly focuses on the random seeds.

	All (N=13)	Meta-Learning (N=14)	Fine-Tuning (N=35)	In-Context Learning (N=96)	Prompt-Based Learning (N=17)	PEFT (N=10)
Label Selection	1	3	4	3	1	
Data Split	2	1	4	3	1	
Task Choice	1	7	N/A			
Data Choice	1	9	6	51	8	
Random Seed	7	3	20	7	3	2
Model Initialisation	9	1	13	N/A	2	4
Order of Data	8	2	11	37	6	2
Implementation	5	1				
Number of Samples		3	9	17	2	2
Hyperparameters	3	1	4	1		
Augmentation	3		1			1
Noise	1		2	1		1
Prompt	N/A	N/A	N/A	46	11	5
Model Randomness	3	1	3	3	2	1
Approach Specific	N/A	1	5	3		

- (2) *Determine* - an effort to discover the origin of the randomness (different from the *randomness factors*) is made by analysing the behaviour of the specific machine learning approach or model in the presence of randomness. For example, to determine why changing the order of samples causes a significant change in performance, the experiments that observe specific problems (e.g., label bias) are designed and analysed.
- (3) *Mitigate* - mitigation strategies for reducing the effects of randomness are proposed, used and their effectiveness is determined. For example, it is proposed to use more deterministic data sampling based on heuristics instead of random sampling.
- (4) *Benchmark* - a benchmark that is specifically designed to take into consideration the effects of randomness is proposed. For example, a benchmark that performs multiple train and test splits and evaluates on all of them.
- (5) *Evaluate* - a more sophisticated, usually statistical, framework for evaluating the experiments investigating the effects of randomness, or determining the effectiveness of mitigation strategies, is proposed. For example, instead of calculating a simple average value, it is proposed to estimate the distribution of results and evaluate the difference using statistical tests.
- (6) *Report* - it is proposed to report more than a single value when running and evaluating experiments and comparing approaches to better deal with the effects of randomness. Not necessarily dependent on experiments for addressing the effects of randomness, but also recommendations on how to improve the reporting of results overall when the randomness has a strong effect on performance. For example, instead of reporting a single average value, the minimum and maximum values over multiple runs should be reported as well.
- (7) *Compare* - better comparison strategies between different approaches in the presence of significant effects of randomness are proposed. For example, instead of comparing single average values, it is proposed to compare distributions of results using statistical tests.

We aggregate the closely related tasks together and organise the survey based on these aggregated groups. Namely, we aggregate together the *investigate* and *evaluate* tasks, as the evaluation strategies are designed specifically for the investigation experiments. In addition, we aggregate the *benchmark*, *compare* and *report* values, as they are closely related.

As the different papers can focus on different, often disjointed, tasks for addressing the effects of randomness, this property does not necessarily provide a distinct division of the papers. For example, the most common combination is to perform *investigation*, *determination* and *mitigation* of a specific *randomness factor* in a single paper. In such cases, the paper is mentioned in this survey multiple times, each time focusing only on the parts relevant to the specific task.

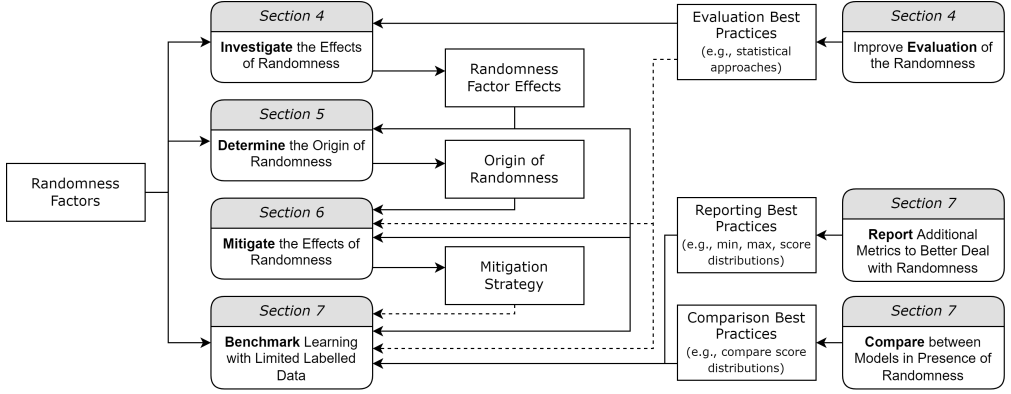


Fig. 4. Different tasks for addressing the effects of randomness, along with their inputs, outputs and the relations between them. The dashed lines represent relations between the tasks that currently do not exist but need to be considered in the optimal state of addressing the effects of randomness, e.g., using more sophisticated evaluation when determining the effectiveness of mitigation strategies. The related tasks are grouped, e.g., *evaluate* task near the *investigate* task.

3.2 Randomness Factors

The *randomness factors* property specifies what *randomness factors* are addressed in the paper in some way (i.e., their effects are investigated or mitigated), but not *factors* that are just mentioned (e.g., mention that other papers found the factor to be unimportant). Papers can focus on multiple *factors* at the same time.

Multiple *randomness factors* have already been identified and described in the existing works (e.g., [44, 110]). Since the used terminology differs across the works, for the purpose of this survey, we established the grounds (in terms of common terminology and understanding) by collecting these *randomness factors* and unifying their denomination and definition.

The identified *randomness factors* can be categorised into three distinct groups: 1) *factors* related to the input data; 2) *factors* related to parameters of the model and the overall training process; 3) *factors* related to the underlying implementation of deep learning frameworks and hardware; and 4) *factors* that provide more systematic changes.

As the *randomness factors* related to the input data, we identify the following:

- (1) *Label selection* - specifies what data is considered to be annotated from a larger set of unlabelled data. Represents the real-world scenario with a limited labelled budget, where only specific samples are annotated. Has a potentially significant impact on the underlying data distribution.
- (2) *Data split* - randomness introduced by different splits of data into training, validation and testing set. Mainly impacts the information provided to the model during training.

- (3) *Task choice* - randomness introduced by choosing different tasks for training. This factor is mainly relevant for approaches that perform *few-shot learning*, such as *meta-learning*, or when using a benchmark dataset composed of different tasks.
- (4) *Data choice* - randomness introduced by the sampling of data the model uses for adaptation (*meta-learning*), or as examples in prompts (*in-context learning*).
- (5) *Order of data* - randomness introduced by the random shuffling of data in different training iterations. Mainly relevant for the *in-context learning* (order of labelled samples in the prompt) and *fine-tuning*. For *meta-learning*, this factor represents the order in which the different tasks are presented during training.
- (6) *Noise* - randomness introduced by adding noise to the different samples.
- (7) *Prompt format* - randomness introduced by small changes in the prompt formats, such as replacing words with synonyms or words without semantic meaning.

Similarly, we identify the following *randomness factors* corresponding to the randomness related to model parameters and the overall training process:

- (1) *Parameter initialisation* - randomness introduced by randomly initialising the parameters of the model, such as weights in neural networks. Mainly relevant for *meta-learning* and *language model fine-tuning*, as no initialisation is performed in *in-context learning*.
- (2) *Random seed* - randomness introduced by non-deterministic actions in the training process. Often used as an aggregated *factor* for the parameter initialisation and the order of data *randomness factors*.
- (3) *Model randomness* - non-deterministic operations in the model, such as the use of non-deterministic layers or regularisation.
- (4) *Approach specific parameters* - aggregates the randomness introduced by *factors* specific for only some of the approaches (e.g., sampling strategy or number of iterations).

Randomness introduced by *implementation of deep learning frameworks and hardware* is a special low-level group, which includes *factors* such as scheduling and floating point precision, auto-selection of primitive operations, how parallel processes are used, changes introduced by different implementation of frameworks such as TensorFlow or PyTorch and their different software versions, but also many other *factors* influencing reproducibility more than the stability of results [44, 110]. Since these *randomness factors* cannot be easily addressed, for this survey, we do not specifically distinguish between them and consider them jointly as one group.

Finally, we also consider *factors* that provide more systematic changes and can be viewed as an edge case with regards to randomness. This group includes factors such as the randomness introduced by augmentation, hyperparameters or the number of samples. Besides these *randomness factors*, we also consider *systematic choices* as they can affect the sensitivity to the effects of randomness as well. These are often a result of the experimental setup and include choices such as the datasets and their characteristics (e.g., monolingual vs. multilingual), or the model size

Each *randomness factor* is associated with a set of its *randomness factor configurations* (sometimes referenced in short as *configurations*). A single *configuration* specifies one state from the set of all possible states the *randomness factor* can attain. For example, a single *randomness factor configuration* specifies one possible permutation of the samples (order), a single set of weights used for initialising the whole model, or a single random value to seed the random number generator. Using this definition, the set of *randomness factor configuration* can be infinite for some *factors* (e.g., initialisation weights) or finite, but extremely large (e.g., permutations of order).

The different *randomness factors* may depend on each other, leading to *interactions* between them [25, 104, 147]. An example of such *interactions* is the relation between data choice and order of data *randomness factors*, where the choice of what data we use in training explicitly affects the

order of samples. Therefore, such interactions can lead to negative, often unforeseen consequences, such as prohibiting reliable freezing of *randomness factor configurations* (e.g., it is not possible to have a fixed order of samples when using different choices of data).

3.3 Machine Learning Approach

Similarly, the *machine learning approach* property specifies which group of approaches is addressed in the paper. Besides the machine learning approaches that define the scope of this survey (Section 2.1), we also consider a generic *machine learning*. This is a special value that contains papers that still deal with the effects of randomness, but only as part of overall supervised learning and do not focus on any of the specific machine learning approaches. The main idea, high-level overview and reference to popular approaches and comprehensive surveys of the machine learning approaches we work in this survey are included in the supplementary material (Section B).

4 Investigating the problem of stability

In this task, the effects of randomness on the stability are investigated, their impact is estimated and the significance of this impact is determined. Each paper investigating the effects of randomness can be characterised by three main aspects: 1) *scope* of the investigation; 2) investigation *methodology*; and 3) *results* from the investigation. The *scope* of the investigation includes decisions about what *factors* to investigate or across what models, datasets and modalities the investigation is performed. It also includes more nuanced specifics like what problem statement is considered (simple classification/regression or sequence generation), or from what setting the data comes from (in-distribution vs. out-of-distribution). Although these decisions have minimal impact on the *methodology*, they strongly govern the obtain *results* and their generality.

The investigation *methodology* determines how the investigation is done. It provides answers to different problems and specifics of the investigation, such as how to perform the investigation (multiple fully randomised runs or keeping some parts static), how many investigation runs to use, how to address interactions between *randomness factors* or how to evaluate the different runs. Again, each decision has a significant impact on the *results* of the investigation, mainly how trustworthy and objective they are.

Finally, the *results* of the investigation specify findings and the behaviour that was observed. Interpreting these *results* is important for informing the further tasks for addressing the effects of randomness (i.e., determining the real origin of randomness and mitigating the effects). They can also be used to provide recommendations on the best use of models in practice in low-data regimes.

4.1 Scope of the investigation

The majority of the papers focus on investigating how the randomness affects the language models (41 out of 56) when they are used to deal with limited labelled data. Most of the focus is on *in-context learning* [36, 75, 83, 126, 174, 180] (32) and only a smaller focus is on *fine-tuning of language models* [33, 56, 89, 128, 177, 184, 187] (9) or *prompt-based learning* (4). Parameter-efficient fine-tuning methods are not investigated at all (only recognised in [25]). The focus on investigating overall *machine learning* approaches (mainly their use with few data and evaluation) [16, 17, 32, 110, 120] and the use of *meta-learning* [2, 15, 28, 104, 129, 164] is significantly lower (10 and 6 respectively). Only a fraction of papers consider multiple approaches at the same time [104, 126, 148, 159].

The significant focus on language models also translates to what modality is addressed the most. **The investigation of text data in NLP problems receives the most focus** (46 out of 56) [16, 33, 36, 56, 75, 83, 89, 120, 126, 128, 164, 174, 177, 180, 184, 187], with images being mostly addressed only in *meta-learning* and *machine learning* approaches (but not all of them)

[2, 15, 17, 28, 110, 129], with only 3 papers considering investigation across both image and text data [32, 138, 191].

The popularity and importance of *randomness factors* heavily depend on the approach. The most focus is dedicated to the data order *randomness factor*, being the most popular for machine learning, fine-tuning and in-context learning [36, 83, 108, 174, 180, 191]. The remaining *randomness factors* are investigated only for specific approaches. For in-context and prompt-based learning, the most popular factors are sample choice and prompt format [119, 127, 147, 171] and, to a certain extent, the number of samples. Parameter initialisation and random seed are most popular for general machine learning and fine-tuning approaches [15, 25, 32, 56, 75, 120, 126, 164, 177, 187]. The implementation randomness factors are popular only for the general machine learning papers [15, 45, 110, 138, 190]. The remaining randomness factors, such as the effects of augmentation [16, 138], choice of what data is labelled [16, 28, 56, 75, 104, 129], split of data [16, 32, 104], noise [16], task choice for meta-learning [2, 129], or other approach-specific factors receive only limited attention.

The scope of the investigation has the strongest effect on how general the final results are. Although the majority of the papers focus on **multiple datasets and models** [2, 7, 15, 17, 28, 56, 75, 126, 164, 174, 177, 180, 187, 191], there are studies that consider only a **single dataset or a single model** [16, 33, 89, 110, 128, 129], limiting the generalisability of the results. **Only in specific cases is the investigation specifically designed to observe the effects of systematic choices**, such as the number of samples or model size, and **interactions between randomness factors** [25, 32, 36, 83, 104, 119, 120, 147, 184]. Finally, **only a minority of papers consider an out-of-distribution setting** [3, 28, 89, 164, 167, 187, 191] or **investigation across different languages** [41, 167, 171].

4.2 Investigation Methodology

Overall, **the differences in how the investigation of the effects of different *randomness factors* is done are small across different papers**, making the investigation quite consistent. In all papers, the investigation is performed by changing the *randomness factor configuration* (such as changing the value of the random seed) and observing the changes in the behaviour of the model and its results. In the early papers, this was done by introducing the change randomly into the *configuration*, without any control over the remaining randomness factors [16, 17, 28, 56, 75, 89, 108, 119, 120, 126, 128, 164, 174, 177, 184, 187]. However, the most common strategy is to exert additional control over the investigation by fixing the *configuration* of the non-investigated randomness factors to a specific value while randomly varying the value for the single investigated factor [15, 36, 45, 83, 86, 110, 127, 138, 142, 171, 180]. In specific cases, the fixed investigation is repeated across different fixed values for the non-investigated factors and aggregated over them to better deal with interactions [104, 148]. Finally, another investigation strategy changes the *configuration* for multiple randomness factors at the same time, and their effects are disentangled at the end [3, 25, 33, 147, 149, 191].

The choice of **how many *randomness factors* to investigate at the same time also plays an important role** in the investigation methodology. Many papers consider only a single *randomness factor* [17, 32, 75, 108, 120, 126, 127]. However, such investigation can lead to biased results due to the interactions between different *randomness factors*. Using multiple *randomness factors* in a single investigation allows for a more detailed exploration of the behaviour of models in the presence of different sources of randomness. Most approaches investigate multiple *randomness factors*, although the number can be quite low (2-3) [15, 16, 25, 28, 33, 36, 56, 83, 89, 119, 128, 171, 174, 177, 180, 184, 187, 191]. Comprehensive investigation of a large number of factors is considered only in specific cases [45, 104, 110, 138, 147, 190].

Due to the potentially significant impact of the interactions between *randomness factors* and systematic choices, it is important to define **how to consider these interactions and how to deal with them** to disentangle the effects of different *randomness factors* and investigate *how the systematic choices affect the findings*. Many papers do not use any specific way to distinguish between the effects of different *randomness factors* [2, 17, 28, 32, 75, 89, 120, 126, 129, 164, 187]. As mentioned before, this may lead to the results being biased and can be traced as a direct cause of some of the open problems discussed in Section 8. A simple way to consider the interactions is to use the fixed investigation [15, 36, 83, 110, 177, 180]. However, such an approach is limited in its investigatory strength, due to the randomness still being present in the choice of the fixed *configuration*. A better way is to search through combinations of *configurations* for all the *factors* at the same time [3, 16, 33, 56, 104, 128, 142, 148, 174, 184]. However, **this significantly increases the computational cost of the whole investigation** [16, 32, 190]. In case of systematic choices, the most straightforward approach is to repeat the investigation for different systematic choices, such as the number of samples or model size [25, 85, 104, 119, 149, 171, 191].

Another important aspect to consider is **how to obtain representative results from the investigation**. The representativeness of results mainly depends on how many *randomness factor configurations* are explored, i.e., how many training and evaluation runs are used during the investigation. Many papers choose the number of different *randomness factor configurations* only arbitrarily, often choosing lower number (3, 5, 40, 41, 86, 100, 1000) without any explicit explanation why the given number was chosen [15–17, 36, 75, 89, 110, 120, 126, 164, 177, 187]. Similarly, when considering multiple *randomness factors* and interactions between them, it is common to arbitrarily select some number (usually small) of *configurations* for each factor and then determine the number of runs as all possible combination of these *configurations* [33, 56, 128, 148, 184]. For example, when considering 3 different *randomness factors*, with 5 different *configurations* for each, the final number of runs is determined as 125 (5^3). In specific cases, all possible *configurations* of a single *factor* can be explored [2, 32, 83, 142, 174, 180] or the number can be determined dynamically (e.g., until something happens or a threshold is overcome) [104, 166, 191]. However, this happens only when the number of possible *configurations* is small or the added computation cost is not as significant.

Majority of the papers **evaluate the effects of randomness from multiple runs based on an aggregated value from these runs**, either as a single value of mean, standard deviation (with standard deviation not even being present in some papers), or number of failed runs [17, 28, 32, 36, 56, 83, 89, 110, 126, 129, 164, 177, 180, 184, 187], or as multiple aggregated values, such as average, worst and best performance, or the difference between them (e.g., difference between worst and best run, or difference between worst and average run) [2, 75, 174]. Some papers define new metrics, such as an average disagreement between models [138, 190], a change in performance when models are ensembled [138], a relative gain [4, 190], or a factor importance [104]. Only some papers evaluate the results of the investigation based on the distribution, mostly as a means for comparing different *randomness factors*, models or the behaviour across datasets [120]. Some papers introduce special evaluation strategies based on statistical approaches that can better distinguish between effects of different *factors* and deal with their interactions, and that can better estimate the results distribution from a small number of repeated runs to make the results more representative [15, 16, 33, 128]. Finally, only a minority of papers compare the learned representation instead of the predictions [11, 138, 190].

4.3 Results from the Investigation

Majority of the papers found that the effects of randomness have significant effects on the performance and stability across all approaches for dealing with limited data, leading to differences in performance as high as 90% [2, 15–17, 32, 33, 36, 56, 83, 89, 126, 128, 164, 174, 180, 184].

Based on the randomness, the rankings in benchmarks and results of the comparison can change significantly [7, 86, 191]. Overall, the algorithmic factors were found to be more important than implementation, but only by a small margin [45, 138, 190]. *In-context* learning was found to be significantly sensitive to prompt format [127, 142, 147, 166, 167], choices of samples [3, 167, 183] and order of choices in multi-choice QA [149, 191], leading to average performance difference of 30% and up to 70%. *Fine-tuning* BERT multiple times and picking the best result, it is possible to outperform its more recent variants [33]. The difference in performance between the worst performing and the best performing set of adaptation data in *meta-learning* can be up to 90% [2]. On the other hand, the observed variance in other experiments was found to be lower, usually only a few percent difference (1%-5%), while still being enough to influence the results and comparisons between different models [75, 110, 120, 177]. However, **the effects of randomness are significantly higher in the out-of-distribution setting** [3, 89, 164, 167, 187, 191].

The investigation results indicate that the impact of interactions, systematic choices and experimental setup on the sensitivity to the effects of randomness and findings is significant in many cases [45, 86, 104, 147, 171, 191]. However, there is usually **low consensus** on how much the variance is reduced by different systematic choices. Some papers found that using more labelled samples or more training iterations reduces the stability [2, 3, 25, 56, 104, 119, 183], although with diminishing returns [85, 177]. On the other hand, other papers found that increasing the number of samples does not lead to significant improvement [83, 108, 119, 167, 191], or even having negative effects on the stability [180]. Similarly, some papers found that increasing the model size reduces the variance [171, 184], but others have found that it either has no effect or even worsens the situation [25, 36]. Optimising *configuration* of specific randomness factors leads to lower sensitivity to other factors, such as using sophisticated sample selection reducing sensitivity to order [2, 174, 183], although it varies across experimental setups [2, 171] (e.g., using 16 random samples is equivalent to using 1 high-quality sample). Keeping the *randomness factor configuration* fixed during the whole use of the model, such as using a single set of samples or defined order based on heuristics, can help reduce the variance [129], but not in all the cases [15, 174]. Finally, accounting for and investigating as many sources of randomness as possible helps with reducing variance at the cost of computational cost [16].

The behaviour of *randomness factors* is affected by the choice of models and datasets. **The inconsistency in results is present here as well.** Some papers indicate that there is different behaviour of the *randomness factors* across different model architectures [10, 17, 32, 83, 86, 120, 149, 171, 184] or datasets [10, 86, 164, 167, 171, 180, 187], with smaller datasets showing more variance [33]. However, many of the results from investigation indicate that the behaviour of *randomness factors* is consistent across all the model architectures, sizes and datasets [2, 15, 16, 36, 56, 75, 89, 126–128, 148, 174, 177].

The investigation of the effects of randomness often just determines that there is a sensitivity to a specific factor, but **without providing any additional analysis of how much it affects the models or what randomness factors are most important** [28, 75, 89, 126, 128, 129, 164, 167, 174, 187]. However, **the most interesting use for the results of the investigation is to provide more in-depth analysis.** For example, comparing between different randomness factors, datasets and models [104]. Although both choice and order are significant contributors, some results indicate that choice is more important as with good samples the sensitivity to order is significantly lower [3, 85, 104]. The prompt format can even affect the sensitivity to the chosen samples [171]. Other results indicate that random seeds, especially the ones used when *fine-tuning*, are more significant than just the effects of implementation and hardware [110, 138, 184, 190], while others find the effects implementation and hardware to contribute more significantly [15]. The results from the investigation of almost all randomness factors at the same time indicate that

the splitting of data has the most impact on the performance [16]. Finally, **analysing the results can better inform other choices later in the process**. For example, estimating the benefit of prompt formats based on closeness (similar to hyperparameter tuning) is not possible, as even the best and worst performing prompts can be close to each other [166]. In addition, the quality of the prompt is not dependent on how well the instructions are written, as even misleading or irrelevant instructions can lead to good results (e.g., better than the correct prompt) [148].

4.4 Overview of Findings: Investigating Effects of Randomness

The effects of randomness are **not investigated evenly across all groups of approaches for dealing with limited labelled data or modalities**. The investigation is more skewed towards randomness in language models, which goes hand in hand with the focus on natural language processing. However, the findings from *meta-learning* and overall *machine learning* approaches indicate that this mostly stems from the popularity of language model research and not from the impact the randomness has on the investigated approaches. **The effects of randomness are consistently significant across all approaches** for dealing with limited labelled data.

In addition, the **focus of the investigation is mostly on a few popular randomness factors**, such as random seeds, model initialisation, sample choice and order, which may cause lower instability. Even though other *randomness factors* receive minimal attention, it has no correlation with their importance. For example, **data split is considered the most important factor**, while being investigated only in a few papers. Another mostly ignored factors are the non-deterministic implementation and hardware, and the choice of which samples are labelled.

The **investigation strategy design is agnostic from data and the machine learning approach**, making its use consistent across all papers with only small modifications. However, these small modifications have a significant impact on the overall investigation, **often making it impossible to compare findings from different papers**. The number of training and evaluation runs are **chosen rather arbitrarily**, ranging from ones to thousands, **without any explanations behind the choice**. Almost no paper introduces heuristics for determining how many repeated runs are needed to estimate the underlying distribution of results. The only heuristic used is to explore all the possible *randomness factor configurations* when their number is small, such as the permutations when the number of used samples is small (< 10).

The effects of randomness are usually **investigated only across a small number of randomness factors at one time**. The number of considered *factors* is usually 1-3, while other *factors* may be mentioned, but not investigated. Only papers dealing with general machine learning cover a broader and more representative set of *factors* at the same time. However, it is problematic to investigate multiple *factors* at the same time due to the computational cost this introduces - the number of repeated runs raises exponentially with a number of *randomness factors* [16, 17, 32].

Even when evaluating multiple *randomness factors* at the same time, **the interactions between them or the impact of systematic choices are not taken into consideration**. Many papers use the fixed *randomness factor configuration* in the investigation, but it is still affected by the randomness. Only some papers use combinations of all *randomness factor combinations*. However, this is always **done only for factors that are considered in the paper** and not for all of them. In addition, no paper applies mitigation strategies to other *randomness factors*.

No analysis of the importance of the *randomness factors* is performed when investigating a single *factor*. Similarly, when investigating multiple *factors* at the same time, **almost no comparison is performed between these factors**. This causes the investigation results to be limited in their extent as the problem is only recognised, but no further conclusions can be drawn from it. It is not possible to answer questions like: “Are the effects of the *factor* significant?” or “Which *factor* is more important and should be focused on first?”.

The effects of randomness are evaluated in a multitude of ways. The most common approach is to use a simple aggregation over the runs (e.g., mean or standard deviation). However, as the experiments are performed on a limited number of training runs, the single aggregated value may not be representative of the underlying distribution. **A better approach for evaluating the effects is the use of statistical approaches** that allow to quantify the uncertainty still present in the models when only limited runs are used, or to determine the expected behaviour of the models in limit, while only having results from a significantly smaller number of runs. The evaluation is further complicated by the interactions between *randomness factors* that can muddle the results and need to be taken into consideration using statistical approaches.

Focus on the more specialised settings, such as out-of-distribution, multilingual, or using parameter-efficient fine-tuning methods, is severely limited. However, few works that study these settings found the effects of randomness to be more significant in such settings, especially for out-of-distribution tasks or parameter-efficient methods.

A significant drawback regarding the results from the investigation is that **no obvious consensus exists regarding the behaviour of different *randomness factors* across different models and datasets**. Many papers investigating the same setting come to contradictory findings about almost all parts of the investigation. Some papers observe a decrease in variance when the size of the model increases, while others found no obvious connection between the size of the model and the amount of observed variance (e.g., increasing size leads to an increase in variance). When it comes to datasets, some papers observed consistent results, with the same *factor* being most significant across different datasets, while others observed different *factors* being important on different datasets. This may be a result of disregarding interactions or using a sub-optimal experimental setup, as some papers observed both interactions and systematic choices have significant effects on the findings of investigation (e.g., the effect of one factor may be misattributed to another factor).

The most consistent finding from the investigation is that **increasing the number of samples reduces variance**. However, the decrease in stability has diminishing returns, where from a certain point, increasing available data has no more significant impact. This may also be due to the quality of samples used, as some papers found that 1 high-quality sample provides a reduction in sensitivity equivalent to 16 randomly selected samples.

5 Determining the Origin of Randomness

In this task, the origin of randomness is determined, which represents a specific underlying problem that only manifests through the *randomness factors*. The design of experiments is similar to the one when investigating the effects of randomness, but with a specific hypothesis that is formed and answered either empirically or using theoretical understanding of deep learning models (e.g., whether the label which is used as the last example in *in-context learning* has an impact on the prediction, which would indicate recency bias). Out of 161 core papers covered by this survey, only a small fraction of 27 papers focus on this task, even though it is important for the mitigation of the sensitivity to randomness.

The design of experiments and further explanation of the behaviour is closely tied to the used *machine learning approach* and the *randomness factor* that was observed to cause the randomness.

5.1 Origins of Randomness in Meta-Learning

The focus on origins of randomness in *meta-learning* is limited only to the adaptation data sampling *randomness factor* [2]. The behaviour of different *meta-learning* approaches is analysed using a combination of all the possible combinations of the adaptation samples from image datasets. The goal of the analysis is to determine the characteristics of the worst-case samples, i.e., the samples that cause the worst performance when used for adaptation. It is determined that these samples are

not artefacts, as they are always correctly labelled, appear representative of the respective classes and are not visually different from other images. Based on this observation, it is speculated that the problem is due to the adversarial nature of the images. However, it is later disproved empirically as adversarial training does not reduce the effects of randomness, leaving the real reason why *meta-learning* is sensitive to the choice of support samples open.

5.2 Origins of Randomness in Language Model Fine-Tuning, Prompt-Based Learning and Parameter-Efficient Fine-Tuning

As the main sensitivity and origins of randomness for *language model fine-tuning*, general *machine learning*, *prompt-based learning* and *parameter-efficient fine-tuning* are related to each other, we describe them together. The origin is analysed mainly using the random seed *randomness factor* (overall across 14 papers), which affects the initialisation, order of samples and model randomness. Many different origins of randomness are identified. The experiments in this part are the most diverse (e.g., models, languages, factors) but often lead to similar origins of randomness.

Small availability of data that causes overfitting and catastrophic forgetting is identified as being one origin of randomness when fine-tuning language models [53, 177]. This can cause the models to rely on shallow lexical features, even though the models were trained on abundant, feature-rich data beforehand [177]. Another possible origin may be the bad fine-tuning start point in the transferred model, mainly in the higher layers (not only the newly initialised ones but also the pre-trained ones) that are sensitive to perturbations in the input.

However, catastrophic forgetting and small size of training data are disputed as causes by later papers [94, 152, 173]. Instead, it is only a symptom of another problem [94]. Looking at the failed fine-tuning runs (i.e., the runs that perform worse than random chance), it is identified that the language models suffer from vanishing gradients. As the cause of the vanishing gradients, a sub-optimal experimental setup is identified, such as using an optimiser without bias correction or a small number of training iterations [94, 173]. An additional cause of vanishing gradients is the sub-optimal initialisation point for the upper layers, which are specialised and cause significant variance when transferred [173].

Again, all these hypotheses are disputed [152] and instead under-specification of the optimisation problem is proposed as the origin of the variance [30, 89, 152], especially in out-of-distribution datasets [89, 152]. There exist multiple good solutions in the generalisation errors on the source dataset, but with wildly different generalisation errors on the target datasets, with the optimal solution lying in a non-flat region where even a small perturbation leads to large difference in error [89, 152]. This makes the identically trained models encode different inductive biases [30] and makes the choice of minima rather arbitrary, being easily affected by small changes in initialisation and order of examples [89]. Under-specification was also identified as the main origin of randomness for prompt-based learning [60] and parameter-efficient fine-tuning (mainly prompt-tuning) [26], which causes the loss surface to be remarkably steep and a small change in input causes a massive change in the loss space.

Furthermore, it is postulated that the large variance in performance is due to the strong inter-example correlations in datasets [187]. When the prediction of a single sample changes, this correlation causes simultaneous change in a large proportion of the predictions in the dataset. This was observed when analysing the training trajectory in an experiment, where the model checkpoints at the same step of training from different initialisations had significant differences in performances. The individual predictions of each sample were highly unstable and changed constantly, which caused large fluctuations. However, no explanation for the unstable behaviour in individual samples is provided, besides being an inevitable consequence of current datasets [187].

Finally, it is claimed that randomness is the inevitable outcome of the stochastic training and a basic amount of instability will be always observed, equivalent to changing 1 bit in the weights [138].

5.3 Origins of Randomness in Prompting and In-Context Learning

The origins of randomness in *in-context learning* are analysed across three *randomness factors* (overall 13 papers): 1) order of data (both samples and answer choices in multi-choice question answering) [83, 180]; 2) choice of samples [75, 174]; and 3) prompt format.

The significant effects of different data orders, sample choices and prompt formats were found to be the consequence of highly unbalanced label distribution in the output [83, 84, 127, 170, 180]. The output distribution of the language models is influenced by the biases present in them [37, 108, 119, 176, 180, 181, 191]. For the order of in-context examples, four biases were identified that cause these effects [37, 119, 176, 180]: 1) majority label bias, where the output is biased towards the more frequent answer in the prompt (e.g., class imbalance in prompt); 2) recency bias, where the model tends to output the labels that appear towards the end of the prompt; and 3) common token bias, where the model is biased to output more common tokens from the pre-training distribution; and 4) domain label bias [37] that is introduced by the domain texts. For order of answer choices, mainly two biases are identified [108, 181, 191]: 1) positional bias, where the model prefers answers in specific positions; and 2) token bias, where model prefers specific answer symbol (such as Option A). The combination of these biases causes the model to have a highly unbalanced label distribution that needs to be calibrated to reduce the variance. In addition, the biases are enhanced by the choice of prompt format and in-context examples [119, 127, 176].

On the other hand, calibrating the distribution to reduce the variance and sensitivity to prompts was found to not be effective, as the bias still persists [83, 191]. This leads to the hypothesis that the biases are not the sole cause of the highly unbalanced label distribution and the origin of randomness should be traced to something different, such as a combination of biases and relationships between the choices and distractor answers [191], or prior in-context examples [170].

The random sampling of examples for the *in-context learning* was determined to be the origin of randomness [75, 174] that can also cause problems observed when using a different order of examples in prompts. Choosing the samples randomly, without taking similarity or other heuristics into consideration, may result in many low-quality demonstrations being selected that may contain spurious correlations and may not provide enough information [75, 174]. When low-quality samples are used, different permutations and calibration of output distribution do not help reduce the variance, which disputes the claims in the previous works [174]. Choosing high-quality samples based on similarity or heuristics reduces the variance and allows different ordering and calibration to work [75, 174]. Therefore, the real origin of randomness, which can also lead to unbalanced distribution and sensitivity to order, is the choice of low-quality, dissimilar samples [75, 174].

Finally, it was observed that sensitivity to prompt format is affected by how frequently the prompt appears in some variation in the data during training [43].

5.4 Overview of Findings: Origins of Randomness

The observed origins of randomness can be summarised as follows:

- **Poor choice of samples** that are used for training, adaptation or in-context learning.
- **Overfitting** that causes **catastrophic forgetting**, and models to focus on shallow features.
- **Under-specification**, where multiple local minima with the same performance are present in training data, which are not consistent with testing data.
- **Highly unbalanced output label distribution** stemming from biases present in language models and incorrect ordering of samples.

- **Optimisation problems** caused by poor choices in the experimental setup, such as not using bias correction, training for a limited number of iterations or poor initialisation and re-use of some neural network layers.
- **Prompt format** mainly caused by biases, such as how often the specific prompt and words in it are encountered during training.

The individual **origins of randomness are closely tied to the *machine learning approach* and the *randomness factor*** investigated. This is mainly the consequence of the experiments that are designed to observe the behaviour of different approaches. Although many different origins are identified, all are tied to a small number of *randomness factors*, specifically a choice of samples, an order of the samples in training, a prompt format and model initialisation.

The popularity of different *machine learning approaches* and *randomness factors* is clearly shown in this task for addressing the effects of randomness. The focus on determining the origins of randomness in *meta-learning* is limited, while for *prompting*, *in-context learning* and their *fine-tuning* it is most popular.

The **first effects of interactions between *randomness factors* are observed**. In many cases, the different origins of randomness identified are later disputed and denoted as a simple consequence of other effects of randomness. For example, it is observed that the choice of samples in *in-context learning* causes problems that are often incorrectly attributed to the ordering of these samples. In addition, different randomness factors were found to enhance the sensitivity, such as prompt format affecting the biases in models that are a source of the sensitivity to data order.

This further shows the **prevailing inconsistency of results** from different papers. The compounding effects between different *randomness factors* and systematic choices cause significant problems when not properly addressed. Another aspect that plays a role is the simple evaluation, as the interactions and systematic choices are clearly shown when a more comprehensive evaluation is utilised.

6 Mitigating the Randomness

In the mitigation task, the effects of randomness are mitigated to reduce the variance as much as possible and to improve the overall performance using mitigation strategies. Mitigation represents the most popular task for addressing the effects of randomness (overall 112 papers). We divide these mitigation strategies into two separate groups: 1) **general** mitigation strategies; and 2) **problem-specific** mitigation strategies.

The general mitigation strategies can be used to mitigate the effects of randomness stemming from any origin of randomness and *randomness factor*. So far, the only way to deal with any origin of randomness is to use ensemble-like approaches (used in overall 13 papers), repeating the training and evaluation multiple times, each time with different *randomness factor configuration*, and then aggregating the predictions (using majority voting or weighted average). The different versions of the ensemble strategy were shown to provide effective mitigation for sample choice [116], model randomness and optimisation [39, 102, 138, 146], order of samples or answer choices [108, 191], the prompt format [6, 55, 142] or data split [96, 182]. In some cases, it is the only available solution for mitigating the randomness, i.e., when the *randomness factor* cannot be easily controlled, such as the effects of *implementation of deep learning frameworks and hardware*. However, it introduces a significant increase to computation costs [16]. Only some papers focus on reducing the cost of the ensemble strategy [102, 138].

On the other hand, the **problem-specific** mitigation strategies are designed to deal with a specific origin of randomness, making them dependent on the *randomness factor* and the origins of

randomness as analysed in Section 5, but also more popular (used in overall 99 papers). In addition, these approaches often provide benefits only for the specific setup (e.g., task type, dataset, or model).

6.1 Poor Choice of Samples

Addressing the sensitivity to the choice of samples is the most popular approach (overall 51 papers). Samples for different approaches are chosen using two approaches: 1) selecting representative samples for the whole process of training and evaluation; and 2) sampling data for the specific training and evaluation iteration. The first approach represents the real-world scenario where there is only a limited labelling budget available that needs to be effectively utilised by choosing samples with as much information as possible. For *in-context learning* this approach is designed to select a single set of examples that perform well for all test samples, instead of selecting separately for each test sample. In the second approach, the samples are selected for a specific training iteration (for fine-tuning) or for a specific test example (for *in-context learning*).

6.1.1 Selecting Representative Samples. When selecting a small number of representative samples (< 100), strategies that work with unlabelled data are often used. This includes the use of a simple clustering [20, 27, 60, 151, 165]. To choose K samples, the available data are clustered based on their similarity [20, 27], diversity [60, 151] or the pseudo-labels generated by a large language model [165] into K clusters. Afterwards, samples are selected from each cluster either as the closest to the centre [20, 27], or based on weighting using uncertainty [165] or KL divergence [60].

Another possibility is to use active learning strategies to select the samples based on uncertainty, diversity or entropy [87, 157], or their combination with other properties such as complexity or quality [77, 88]. Which combination of properties to consider can also be determined dynamically based on the approach used [105]. The submodular functions are often used to optimise multiple properties at the same time [54, 114]. Another possibility to select the samples is to use reinforcement learning, similarly to active learning [174].

A popular strategy is to select samples based on their quality. To determine the quality of samples, one solution is to train a linear regression model that determines the gain in performance when the sample with given characteristics is included in the process [22, 56, 97, 122]. Another possibility is to prompt a large language model to rate each sample based on how it contributes [69, 74].

Another set of strategies selects samples in multiple steps. In each step, a separate sample property is optimised, such as first choosing the most diverse or informative samples and then refining them using their similarity to test sample or their quality [70, 137, 162, 163]. To achieve this, determinantal point processes are often used [162, 163], a large language model is prompted to provide quality/informativeness scores for the samples [70], or a graph-based approach is used [137].

Finally, specific mitigation strategies focus on selecting the samples, their order and the format of the prompt at the same time, using neural bandits [153].

6.1.2 Sampling Data for Training Iteration. When selecting in-context examples for a specific training or evaluation iteration, the most common approach is to use the similarity to test sample [3, 8, 24, 75]. The similarity selection is often combined with other properties, such as text characteristics [154, 167], entropy [1, 106] or margin [189]. Besides using different properties, a better representation can also be used, such as using large language models to transform each sample into how well it represents a certain skill [8].

Another possibility is to use adversarial training on the worst-case adaptation samples (found using greedy search) in *meta-learning* [2]. However, it provides no significant improvement in regards to stability for specific *meta-learning* approaches [2].

6.2 Overfitting and Catastrophic Forgetting

How the overfitting and catastrophic forgetting is mitigated is closely tied to the approach for which it is designed. Therefore, we provide further categorisation based on these approaches.

6.2.1 Overfitting in Meta-Learning. One solution to reduce the overfitting in meta-learning approaches is the augmentation of the support and query data, and the task construction [99]. Applying such augmentation increases the diversity of the samples and tasks, which improves the overall generalisation and reduces the sensitivity to the randomness [99].

Another solution is the use of variance reduction algorithms in the optimisation-based meta-learning [145, 160]. One possibility is to focus on first-order meta-learning algorithms and introduce the variance reduction only to them [145]. The variance reduction term is introduced to the gradient estimator in the task adaptation stage, motivated by the recursive momentum technique in [29]. This term is initialised by averaging the gradients from randomly sampled tasks using the initial parameters and then is further updated using a weighted sum of mini-batch stochastic gradients across all sampled tasks [145]. As this foregoes the benefits of the bi-level meta-learning process, other papers propose to modify the variance reduction term (STORM [29, 66]) for use on bi-level meta-learning optimisation, in combination with a large number of training steps to reduce the amount of storage required [160]. To achieve this, the variance-reduced gradient estimate is computed using stochastic gradients from two successive iterates while evaluating the gradient using the previous two iterates on the current batch of samples [160].

6.2.2 Overfitting in Language Model Fine-Tuning. The augmentation technique, in combination with regularisation, is also proposed when fine-tuning language models [90]. The pre-trained language model is used as a generator for new data samples used for augmentation. To guarantee that samples with discriminative labels are generated, a meta-weighted maximum likelihood objective can be used for the tuning of the generator. As the resulting labels still contain some level of noise, a noise-robust procedure for training the final model is proposed as regularisation using label smoothing and temporal ensembling [90].

Another solution is to look at the disconnection of the problem statement between the pre-training and fine-tuning stages. When the objective or even the input and output distributions are different, the fine-tuning process can often be brittle, especially when only a small amount of data is available [111]. A second pre-training step on data-rich supervised tasks can be used to better prepare the model for the fine-tuning stage, mitigating the brittleness and tendency to overfit on data [111].

The last set of approaches utilises regularisation to prevent overfitting. Injecting a Gaussian-like noise vector into the input, either before or during training, can be used as a stability regularisation, similar to augmentation [53, 150]. Another regularisation technique is to perform fine-tuning updates only on a sub-network instead of the whole network [156, 168, 169]. Instead of random selection, CHILD-TUNING, introduced in [156], uses Fisher Information to determine the importance of parameters based on all samples before training. Afterwards, an unchanged sub-network with the most important parameters is selected for the update. The previous two approaches can be improved using Dynamic Parameter Selection, introduced in [169], which adaptively selects a promising sub-network to perform the fine-tuning steps on. The sub-network is chosen dynamically, i.e., different sub-network is potentially fine-tuned at each training step, based on the importance of the parameters. It greatly reduced the heavy computation costs introduced by CHILD-TUNING [156] that separate the process of fine-tuning and the decision of which sub-network to optimise. In addition, it can better capture the task-relevant parameters across sub-networks of different sizes and bring stable improvement in various parts of the model. However, when faced with limited

resources, the previous approaches cannot select optimal sub-network [133]. As such, the previous approach is modified to represent the weights as a weighted mixup of pre-trained weights and the optimised weights, which is guided and optimised using attention.

6.3 Under-Specification

To deal with under-specification, models with stronger inductive biases that can help distinguish between many local minima should be proposed [89, 187], or the training sets should be better designed [89, 187]. The better sets should represent the phenomena appearing in texts to reduce the probability of local minima not transferring to testing data [89]. In addition, the tasks used for large-scale pre-training should adhere to structure and level of compositionality (i.e., solving a combination of multiple tasks leading to the solution for a different, more complex task), and more diverse datasets in terms of syntax and lexicon should be constructed [187].

Another solution for dealing with under-specification is to utilise an ensemble of models [48, 59]. One possibility to construct the ensemble is to combine predictions from multiple trained models using their uniform average [48]. Another possibility is to use stochastic weight averaging (SWA) as regularisation for the ensembles [59]. Instead of training the model multiple times, SWA creates the ensemble by averaging over different snapshots of model weights. Such an ensemble explores other solutions around the local minima, while also reducing computation cost as the ensembling at weight level needs only a single training run [59]. Although the SWA method consistently achieves better stability of performance than standard fine-tuning when used over different values of random seeds, the increase in stability is only slight [59]. In addition to ensembling, distillation and co-distillation can also be used [48].

To deal with the under-specification in prompt-based learning and prompt-tuning, the best approach is to train over different prompts at the same time [60], creating a kind of ensemble. Another well-performing approach is to train using noise regularisation, adding noise to the input text and embeddings [26].

6.4 Highly Unbalanced Label Distribution

To deal with highly unbalanced label distribution in the *in-context learning* it is proposed to use calibration. One possibility is to use contextual calibration, where the output distribution is first estimated using neutral prompts [37, 74, 180, 186]. The neutral prompts contain only in-context samples with the testing sample having its text replaced with words without semantic meaning, such as “N/A”. Using multiple such prompts, the distribution between labels can be estimated and then corrected using affine transformation [180]. For example, if the model outputs positive labels in more than 50% of cases (e.g., in 65%) in binary classification, the testing example is considered to be positive only when the confidence is higher than this threshold (e.g., more than 65%). Besides words without semantic meaning, other studies observed that sampling from the text corpus [37] or the specific batch [186] leads to better estimation of the distribution.

Instead of using neutral prompts, it is possible to parameterise the distribution using Bayes rule [91]. The language model is not presented with K concatenated samples in a single prompt but instead using one sample at a time. Obtaining output probabilities and multiplying them together can be used to estimate the distribution, while also reducing the dependency on the ordering of samples [91]. Compared to direct distribution estimation, this provides larger improvements in terms of stability when there are fewer training samples, the data is imbalanced, the number of classes is large, or a generalisation to unseen labels is required [91].

However, other studies found that previous calibration strategies do not lead to significant mitigation of the biases [119]. Instead, they propose to do a leave-one-out calibration using the in-context examples, where the probability is estimated by leaving out the specific in-context

examples [119]. Adding noise to the sample parameters was found to be another effective calibration technique [179].

The calibration approaches are used to deal with the sensitivity to the order of choices in multi-choice question answering as well [149, 181, 191]. However, calibration was again found to be sub-optimal, leading only to negligible improvements [191]. The best-performing strategy was ensemble, where the evaluation is run multiple times with different orderings, and then using majority voting [108, 191], or fine-tuning the model further using different choice orders [191].

Apart from calibration, the highly unbalanced distribution is addressed by finding the best ordering of examples [61, 80, 83, 158]. One possible approach to automatically generate good prompt orderings is to use heuristics and entropy [83] or KL divergence [158]. Another possibility is to search over possible permutations of training samples to find the best performing one using genetic algorithm [61]. Instead of greedy-search, curriculum learning can be used as well, to order the in-context examples from easiest to hardest [80].

Finally, a simple solution is to just use longer context and more samples, as it leads to lower sensitivity to order [13].

6.5 Optimisation Problems

The mitigation strategies for dealing with optimisation problems in *language model fine-tuning* are simply specific fixes for the poor choices that introduce the effects of randomness. One of the fixes is to increase the number of iterations considerably and introduce a bias correction and small learning rate to avoid vanishing gradients [94, 173].

Another solution is to modify the way the additional layers are initialised [31]. Instead of initialising the layers randomly, which can lead to bad initialisation, a meta-learning approach can be used to find an initialisation that is amenable to learning using gradient descent [31]. In addition, badly performing initialisation can be detected as early as the second training iteration thanks to the high correlation of performance at the start and the end of the fine-tuning process [33]. Starting many runs with different initialisations and stopping the ones with bad performance in the early stages of fine-tuning, it is possible to search through the whole space of random initialisation to find the best performing ones without any significant increase in computation cost [33].

Yet another solution is to also randomly re-initialise the top layers of the language models, in addition to the classification layer [173]. The number of re-initialised layers depends on the problem statement and the model, with more complex models and datasets benefiting from a larger number of re-initialised layers, but helps only up to a certain point [173].

6.6 Prompt Format

Another popular focus is on mitigating the sensitivity to prompt formats (overall 32 papers). One strategy for dealing with the prompt format sensitivity is automatically selecting good performing prompts [43, 131, 135, 175]. Starting from a seed prompt format, new prompts are automatically generated as a paraphrase using large language models or using back-translation. The generated prompts are ranked based on how well they perform on development set [131], training a separate ranking model [175], or using heuristics such as perplexity [43] or mutual information [135].

Besides generating new prompts, another possibility is to optimise the seed prompt directly. To achieve this, the influential words in the prompt are identified and then iteratively replaced using unmasking [166], starting from the most influential ones. Instead of replacing the words one at a time, it is possible to use a large language model to rate the prompt on a specific mini-batch and suggest changes to it to achieve lower sensitivity [113].

When optimising the prompt, a popular strategy is to add soft prompt embeddings that are then further optimised [139]. However, using soft prompts was observed to be significantly sensitive

to their initialisation. To deal with this sensitivity, a meta-learning approach is utilised to find a good initialisation [50, 101, 115]. Besides meta-learning, the initial prompt can be improved by training it on a few source tasks and then fine-tuning it on the target task [143], or optimising its representation using a separate neural network [78, 118].

A popular approach is to use multiple prompt formats at the same time. Similarly to the previous approaches, the seed prompt is used to generate new ones, e.g., by paraphrasing the hard ones. These new prompts are then used to fine-tune the models on all of them at the same time [24, 67, 124, 185]. Another possibility is to use these prompts in an ensemble, where each prompt is used to obtain a prediction and the final evaluation is done as a majority voting or (weighted) average of the partial predictions [6, 55, 142].

6.7 Overview of Findings: Mitigating Origins of Randomness

Mitigation strategies already receive the majority of the focus from the tasks for addressing the effects of randomness (112 papers). Overall, the main focus of mitigation is dedicated to *prompting* approaches (including *in-context learning* and *prompt-based learning*) (83 papers) and their sensitivity to prompt format, choice of samples and the order of samples and answer choices.

The proposed **mitigation strategies are closely tied to the origin of the randomness**. As they are explicitly designed to deal with the specific origin of randomness and the *randomness factor* it was observed around, the proposed strategies are unusable for other *randomness factors* and origins of randomness. The **best-performing mitigation strategies are ensembling and further model fine-tuning**, shown to provide benefit for many different randomness factors. **Some randomness factors can be mitigated only using the ensemble strategy**, such as the non-deterministic *implementation of deep learning frameworks and hardware*. As such, it is an ideal mitigation strategy, when the introduced computation costs are not a problem. Specific mitigation strategies should be used when only a specific factor needs to be mitigated while minimising the cost of the mitigation.

Even though the mitigation strategies address the specific origin of randomness, **many papers do not provide any kind of extensive analysis of where the randomness originates from**. In some cases, the works reference other papers that found this origin, but **in many cases no explanation is provided as to why the mitigation strategy targets the specific origin of randomness**.

Each mitigation strategy incurs additional computation cost. The most significant increase is introduced by the ensemble mitigation strategy, as it requires a large number of model training and evaluation runs (the mitigation effectiveness depends on the number of models in the ensemble). The expected computation cost is several times higher than in mitigation strategies that require only a single training run. So far, no general mitigation strategy exists that could provide mitigation without high computational costs. In the case of problem-specific mitigation strategies, the computation cost depends on multiple aspects, such as how often it is applied (once before the training run or at each training step), how expensive the approach used (genetic algorithm or simple clustering), or how many repeated runs are needed (e.g., when ensembling is used).

Effectiveness of the mitigation strategies is evaluated in a simple fashion by comparing the performance metric of the model after applying the mitigation strategy. The comparison is usually done on a single aggregated value, often with a low number of repeated runs (sometimes even one), while foregoing the reporting of standard deviation or confidence intervals. In addition, the **interactions between factors and systematic choices are often completely disregarded**.

Such a simple evaluation that disregards the interactions and systematic choices may be the **main reason for inconsistency in results**. Similar to the *investigate* (Section 4) and *determine*

(Section 5) tasks, multiple papers report contradictory findings. The mitigation strategy that works in one paper provides no benefits in another.

Finally, **only a single mitigation strategy is often applied at one time in most papers**. Only a few papers use multiple strategies and compare their benefits.

7 Comparing and Benchmarking: The Impact of Randomness in Low Data Regimes

In this task, the effects of randomness are taken into consideration when comparing and benchmarking different approaches for learning with limited labelled data. As the effects of randomness can have a significant effect when comparing, or when designing benchmarks, modifications to the process of benchmarking and comparing are required. This is recognised by the researchers, as the number of papers focusing on these tasks is growing (overall 33 out of 161 papers). Even though the benchmarks and comparisons between approaches are closely related, the proposed modifications that can deal with the effects of randomness look at the problem from different perspectives.

When simply **comparing** the approaches, the modifications are introduced at the end of training. The design of the modified comparisons is independent of the training process and, in theory, can be applied to any dataset, model, or even benchmarks. The independence from the training process makes the modifications easier to design, but often requires specific assumptions, such as performing multiple runs. The modifications are usually designed to be simple, often taking the form of best practices for the comparisons (e.g., use statistical significance tests on the distribution of results). An important aspect of the comparisons is what metrics are reported.

On the other hand, the **benchmarks** are specifically designed to take the randomness into consideration from the start of training. The whole training process is modified based on desiderata that address the randomness, making benchmarks closely tied to the training process. Even though it makes the benchmarks harder to design, they often provide better explanations and comparative power as no assumptions and constraints are needed.

7.1 Comparing Approaches in Low Data Regimes

Comparisons between different approaches in low data regimes are significantly affected by the effects of randomness [103]. Changing the *randomness factor configuration*, the ranking of models can vary significantly (as shown in Section 4.3) [7, 32, 33, 46, 92, 120]. In addition, when using only a small number of runs, there is a high possibility of cherry-picking results that confirm our hypothesis [45].

To deal with the effects of randomness, many papers propose to compare on multiple evaluation runs to achieve unbiased comparisons [16, 17, 32, 35, 120, 126, 128]. Especially for prompt formats, where it was observed that one format is not ideal for different models [127, 136, 142], it is important to compare models across multiple runs [7, 46, 92, 95, 127, 136, 142], keeping as much of the setup (e.g., size of models), using either the distribution of results or the best performance.

Comparing based on score distribution across multiple runs allows us to perform comparisons using statistical significance testing, which is a simple but powerful modification for the comparisons [16, 17, 32, 35, 86, 120, 128]. Drawing conclusions from such comparisons reduces the bias and the risk of rejecting promising approaches or falsely accepting weaker ones [45, 120, 126].

However, comparing between score distributions may not be as straightforward in all cases [16, 35, 128]. To provide truly unbiased comparisons, as many sources of randomness (i.e., *randomness factors*) as possible should be randomised [16, 120]. This not only leads to a better evaluation of the variance associated with each source of randomness but the error of the expected performance is also reduced. However, such comparisons start to become problematic when the evaluated models are not under our control (i.e., behind API) and the only controllable variance is the data we use. In addition, accounting for all sources of randomness is often not feasible due to the computational

cost of running training multiple times [16]. To reduce the cost, a statistical bootstrap approach can be used to estimate the behaviour pm different settings from lower number of runs [112, 128].

Besides simple metrics, multiple papers propose to report additional ones that can allow for more unbiased comparisons using statistical tests. Such metrics include standard deviation, maximum, minimum, mean and confidence intervals [18, 126], the score distributions and individual episode results from multiple runs, data splits, prompt formats or even across different datasets [18, 32, 35, 92, 96, 127, 182]. Last, but not least, reporting how the ranking of the model changes when introduces a small change was found to be important for large language models[7, 92].

7.2 Benchmarks Taking Randomness Into Consideration

Multiple benchmarks designed for evaluating *few-shot learning* are starting to take the problem of randomness and its effects on the stability into consideration [18, 36, 96, 130, 164, 182]. However, the extent to which the randomness is addressed is different across the benchmarks.

In the most simple case, only some of the properties related to stability and/or randomness are included in the design of the benchmark [18, 96, 130, 164]. This includes focus on variable number of shots and classes in *few-shot learning* to better simulate real-world scenarios and to test robustness of models to these systematic changes [18] or comparing the performance from multiple runs with the human performance [96] to allow for paired statistical tests [18]. Another modification is providing means to evaluate both in-distribution, as well as out-of-distribution [130], or looking at the behaviour of models in cross-task settings [164].

Later benchmarks also explicitly encourage the evaluation on multiple training runs [96, 164, 182]. Such benchmarks mostly provide multiple splits of data for training and evaluation [96, 164, 182]. To explore the behaviour in cross-task settings, the splitting is also done on the level of tasks, with some tasks being used only for testing purposes [164].

Finally, some benchmarks are specifically designed to determine the sensitivity of the approaches to the effects of small, random changes [4, 36, 73, 125, 159, 188]. One benchmark focuses on simpler tasks that can be quickly and automatically evaluated, and are trivial for humans (i.e., tasks that elementary school students should perform perfectly at) [36]. For such tasks, the sensitivity to effects of randomness can be more easily and directly estimated and a sensitivity score can be estimated that is used to explicitly penalise models for their brittleness [36]. For more complex tasks, the evaluation is performed across different small perturbations (mainly to prompt format) and a specific sensitivity metric is designed, such as the spread of performance over these changes or the mean relative gain of the model [4, 73, 125, 159, 188].

7.3 Overview of Findings: Comparisons and Benchmarks

The comparisons are specifically designed to deal with the significant effects of randomness when dealing with limited labels. **Changing the randomness factor configuration, such as using different prompt formats, leads to significant changes in rankings.** Even when comparing across multiple training runs, it was observed that **using a single aggregated value leads to biased comparisons.** The most common modification is to compare between result distributions using statistical tests instead. Although this leads to more unbiased comparisons, there are still problems present. In the case of approaches that use prompts, the best solution is to optimise the prompt format for each model (as the specific format is ideal only for specific models), reporting either the spread across different formats or the best performance from the different runs.

Such comparisons are sensitive to the number of training runs, as the **distribution of results can be objectively determined only with a larger number of runs.** However, running multiple training runs was found to be infeasible in many papers. To deal with this problem, approaches

that approximate the distribution without the need for many training runs were proposed recently, such as statistical bootstrapping.

The modified comparisons are utilised only for comparing between different approaches across datasets and tasks. However, it is **important to also compare between effects of randomness factors**, while taking the interactions into consideration. Even though not specifically designed for it, the modifications should be applicable to this problem as well.

When it comes to benchmarks, the randomness is only starting to be taken into consideration, leading to very simple ways to address it. **Most benchmarks simply evaluate across multiple splits, training runs or prompt formats, while keeping the number of runs quite small**, or simply explore some properties related to stability. In addition, even though multiple training runs are performed, many benchmarks compare the models only on a single aggregated value. **Only few benchmarks explicitly measures the sensitivity of models to effects of randomness** and use it to penalise the final performance metric.

8 Challenges, Open Problems and Next Steps

Based on the findings from the comprehensive literature overview presented in the previous chapters, we identify seven challenges and open problems for addressing the effects of randomness on the stability of learning with limited labelled data.

8.1 Focus on Small Fragment of Randomness

The scope of investigating the effects of randomness, determining its origin and the mitigation of these effects is limited in its extent. The problem can be observed on multiple levels:

- *Approaches and modality* - the focus is mostly on the approaches utilising language models (*fine-tuning* and *in-context learning*), while focus on other machine learning approaches, such as *meta-learning* or even using *parameter-efficient fine-tuning* methods, is severely limited in its extent. The majority of the focus is on natural language processing, as it is the modality used in language models, limiting the analysis of behaviour on other modalities. Therefore the effects of randomness are addressed only on a specific sub-part of learning with limited labelled data, even though the randomness also has significant effects on other machine learning approaches (investigated in this survey).
- *Randomness factors* - many *factors* receive limited attention, such as the effects from splitting data, choosing what data is labelled, model randomness, or the impact of non-deterministic implementation and hardware. Instead, the focus is on *factors* in later stages of the training process, such as data order, prompt format or choice of data. This limits the analysis of randomness to a few popular *randomness factors*. The under-researched *factors* may still be important, but there is no consensus due to the limited focus. Some papers claim the effects of these *factors* to be insignificant [17], while others consider them to be the most significant contributors of variance [16, 96, 182].
- *Number of randomness factors* - the experiments are designed to jointly analyse only a small number of *factors* (1-3). Such experiments compound the limited focus on some *factors* and may leave the results of the investigation vulnerable to unaddressed effects of other sources of randomness that may undesirably skew the results.
- *Settings and tasks* - the experiments mainly consider in-distribution and monolingual settings for classification. Only a few papers deal with more specialised settings, such as out-of-distribution, multilingual, use of parameter-efficient fine-tuning methods, or more specialised tasks (summarisation, multi-choice question answering).

We argue that the scope should be widened across all tasks for addressing the effects of randomness, to better explore the behaviour of different *machine learning approaches* across different modalities and *randomness factors*. A simple extension would be to perform multiple experiments, each for different *randomness factor* across all the different approaches the factor is relevant for. A more in-depth analysis would be to analyse the effects of randomness across all *factors* in a single experiment. However, increasing the number of *randomness factors* that are considered in the investigation significantly increases the computation cost [16]. To overcome this problem, more sophisticated experiments should be designed.

Focusing on multiple randomness factors at the same time, across multiple approaches is only starting to appear, even though it allows us to draw more interesting conclusions and better focus on the most important factors [45, 103, 104].

8.2 Limited Analysis and Explanation of Importance for the Effects of Randomness Factors

A significant part missing when investigating the effects of randomness is a thorough analysis of the importance of the different *randomness factors*. In many cases, the importance of *randomness factor* may be obvious, such as when it causes a large variance in results ($> 10\%$). When the observed variance is smaller, a more detailed analysis and explanation is required. The effects measured as 1% standard deviation have different significance when the performance is in high 90% than when it is in low 20%. However, almost all papers simply present a single aggregated value without further details or analysis.

An important missing analysis is the comparison of importance for different *randomness factors*. As it is common to work with a limited computational budget in real-world scenarios (including a growing effort of green AI that aims to limit its CO2 footprint), identifying the most important *factors* would allow us to focus on them and help better allocate the resources for dealing with the effects of randomness. However only few papers compare between *factors*, such as [16, 17, 56, 104, 110, 177, 180, 184].

Finally, there is only limited analysis of how the sensitivity to the effects of different randomness factors changes across different approaches, datasets and models. Specific studies have already found that the importance of *randomness factors* is not consistent across different approaches and modalities. Similarly, the analysis of how different systematic choices, such as the number of labelled samples or size of the dataset, affect the importance of randomness factors is largely missing, although some studies have started to focus on this problem.

8.3 Disregarding Interactions Between Randomness Factors

Ignoring interactions between different *randomness factors* can lead to biased results. We can observe that many papers often find contradictory findings across all the tasks for addressing the effects of randomness. The effects of interactions are most notable when evaluating the mitigation strategies, as this task has the least developed evaluation strategies, often using simple evaluation on a single *randomness factor*. We argue that the interactions between *randomness factors* are one of the primary contributors to the inconsistency in findings, as was already observed for *in-context learning* (using optimised prompt format and high-quality samples making the sensitivity to data order disappear).

The interactions should be taken into consideration in every task for addressing the effects of randomness, whether it is investigation of the effects, identifying the origin of randomness, proposing mitigation strategies, or even comparing and benchmarking different *randomness factors* and areas of machine learning.

One possibility is to include all randomness factors in the process, such as exploring through all the combinations of the *randomness factor configurations* for each *randomness factor*. However, this incurs significant computational cost due to the number of runs growing exponentially with the number of *randomness factors* and so may be infeasible. Some papers already do this, but only for a set of selected *randomness factors* with low number of *randomness factor configurations* [56, 177, 180].

Another possibility is to mitigate the effects of *randomness factors* that are not currently investigated using mitigation strategies. However, such approaches are only starting to appear [104, 148].

Taking interactions into consideration also has an impact on the evaluation of the effects of randomness, estimation of the mitigation strategy effectiveness and comparisons between approaches and *randomness factors*. Each such tasks need to be explicitly modified to take the interactions into consideration. Specifically, the effects of different *randomness factors* need to be disentangled. Statistical approaches, designed in some papers, that can disentangle these effects already show promising results for dealing with this problem [16, 33].

8.4 Oversimplified Techniques and Metrics Employed to Evaluate the Effects of Randomness

The majority of the papers evaluate the experiments when addressing the effects of randomness in a basic manner, often using a single aggregated value based on the model prediction. As the number of runs is chosen arbitrarily, the estimation for the real distribution of the results may be skewed. Therefore, the evaluation is open to the still present effect of randomness which can lead to biased results.

More sophisticated evaluation strategies (based on statistical approaches) that can account for the still present uncertainty in the limited runs, or the ones that can better estimate the real distribution of results without the need for a large number of samples (training runs), need to be designed to deal with this problem. Currently, only a few works determine the effects of randomness based on such statistical approaches [16, 33, 128]. At the same time, these evaluation strategies need to be used throughout all the tasks for addressing the effects of randomness. Even though some more sophisticated approaches already exist, many tasks (such as evaluating the effectiveness of mitigation strategies) still use a basic evaluation based on aggregated value, often from a single run.

Another avenue for improvement is to design metrics that evaluate the effects of randomness based on the learned representation instead of using the metric that observes changes in model prediction and use them as surrogate for estimating the effects.

8.5 Inconsistency in Findings from Different Tasks for Dealing with Effects of Randomness

Our in-depth analysis of this survey revealed that there are many contradictory findings across papers. Even though the effects of randomness are mostly investigated using a consistent approach (i.e., running multiple training and evaluation runs), there are aspects of the investigation methodology that differ across papers. In some cases, it is observed that the *randomness factors* behave in a similar fashion across datasets, modalities and architectures, while in other cases their behaviour is inconsistent [10, 14, 178].

We believe that these inconsistencies are partially a result of disregarding the interactions between *randomness factors* (Section 8.3) as well as using oversimplified techniques and metrics to evaluate effects of randomness (Section 8.4). As such, investigation strategies that can handle the interactions or statistical methods that can better disentangle and estimate the true effects should be used.

In addition, observed inconsistencies may be a result of insufficient evaluation across a low number of samples and runs; or a result of systematic choices. Following, we describe in more detail such sub-optimal and inconsistent decisions in the investigation and evaluation methodologies and how to properly address them to achieve more consistent findings:

- **Low number of test samples and investigation runs.** The effects of randomness are often evaluated using a low number of test samples (100 - 1000) and a low number of repeated runs (10). This often leads to biased results due to the non-representative distribution of the results that is affected by randomness and limits reproducibility. Following the long-used practices in machine learning, both the number of test samples and repeated runs should be significantly higher, especially when dealing with limited labelled data that only strengthens the problems with randomness.
- **Systematic choices.** Although the experimental setup plays a significant role, it is the most inconsistent aspect across studies, such as using a different number of examples, models, model sizes, or hyperparameter setup. As such, the findings cannot be easily compared with each other. Many solutions for this exist, such as comprehensive reporting for the experimental setup, keeping the setup as similar as possible to previous studies, or designing and using an investigation framework that handles this.

Following the analogy of benchmarks, we imagine that many different investigation frameworks can be proposed, with their designs reflecting some specific aspect deemed to be important (e.g., handling interactions, comparing effects of different *randomness factors*, computational cost). However, so far no such explicit, fully fleshed out framework exists, only starting points of ones, mainly in papers introducing better evaluation strategies [15, 16, 33, 128].

8.6 Absence of Effective and Efficient Mitigation Strategies for Multiple Randomness Factors

Many partial solutions for mitigating the effects of different *randomness factors* exist. Such mitigation strategies are explicitly designed to deal with specific origin of randomness and are always used one at a time. However, many *randomness factors* have no specific mitigation strategy designed for them. In some cases, a specific mitigation strategy may not even be possible, such as when dealing with non-deterministic implementation and hardware.

The only mitigation strategy that is usable for all *randomness factors* is the *ensemble* and, in some cases, further fine-tuning of the models. Although these strategies can effectively mitigate the effects of randomness, they incur heavy computational costs, which makes them infeasible in combination with complex models. In addition, they often disregard interactions between randomness factors.

Mitigation strategies have the potential to have the highest impact on all tasks for dealing with the effects of randomness and on the use of learning with limited labelled data in practice. Having an effective and efficient mitigation strategy that can take the interactions into consideration will positively impact all open problems mentioned in this section. Therefore, it is important to focus on designing a mitigation strategy that can address the problem as a whole. The focus should be on effectively combining different *randomness factor* specific mitigation strategies. Besides taking interactions into consideration and achieving a significant reduction in randomness, a specific focus should also be given to reducing the additional computation cost that is incurred. This will be especially important for the *ensemble* strategy, as we see it being used for *randomness factors* without specific strategies designed for them. Some papers, also from related domains, focus on reducing the cost of the ensemble strategy [21, 72, 102, 138].

8.7 Limited Consideration For Randomness in Comparisons and Benchmarks

Even though many modifications are proposed that can deal with the effects of randomness, they are applied only in specific cases of comparisons between approaches, but not as a part of tasks for addressing the effects of randomness. The results from experiments investigating the *randomness factors* and especially determining the effectiveness of mitigation strategies are compared in a simple manner based only on aggregated value. When comparing mitigation strategies, it is common to even use a single run, which may obscure their benefit [102, 105]. Similarly to comparing approaches, the experiments investigating randomness are still sensitive to the effects of randomness as not all *randomness factors* are considered. More sophisticated comparison strategies should be applied when comparing between different *randomness factors*.

In addition, the benchmarks are only starting to consider randomness in their design, even though they are significantly affected by the randomness. Even when randomness is considered in the design, it is usually in a simple manner, such as looking at simple characteristics, providing multiple splits to train on, using different prompts, or repeating the evaluation multiple times with small changes [7, 86, 96, 182]. Even this little consideration is a step in a good direction, but we need benchmarks that explicitly take variance into consideration and use it to penalise approaches that are brittle in the presence of randomness, such as in [36]. Finally, some *randomness factors* are not even considered in benchmarks, such as model initialisation or the choice of data. Therefore, we still have a long way to go to have benchmarks that take into consideration all the important *randomness factors*.

9 Conclusion

In this paper, we have filled the gap of a missing, comprehensive literature survey focused on the existing research works addressing the effects of randomness that negatively impact the stability in real-world settings when labelled data are lacking. We summarised the current approaches that address the effects of randomness by investigating these effects, determining their origin and designing mitigation strategies for different *randomness factors*. In addition, we analysed a specific group of papers that aim to take randomness into consideration when comparing and benchmarking different machine-learning models. Based on the findings from the comprehensive analysis and synthesis of existing works, we identified a list of seven major challenges and open problems when dealing with the effects of randomness and provided some insights into how to handle them.

Although we focused on a comprehensive evaluation, the area of addressing randomness is currently quickly growing and has attracted a lot of attention recently, as many researchers recognised the sensitivity of different approaches to the effects of randomness. Since 2023, the number of papers addressing the sensitivity to the effects of randomness has almost tripled, with most of them focusing on in-context learning and mitigation. As such, we expect that new research works will be continuously released, the area will be evolving and hopefully the challenges and open problems will be addressed.

This survey offers a long-term value as the identified challenges and open problems are integral to the area and have not been addressed sufficiently even with the recent increase in focus. We hope that this survey will help researchers more effectively understand the negative effects of randomness, the tasks performed when dealing with them, grasp its core challenges and better focus the attention to addressing the randomness and the open problems so that the field can be advanced. We expect that it will encourage researchers to focus on the underrepresented *randomness factors*, their interactions, systematic choices and settings so that the effects of randomness and the behaviour of models when they are present can be understood in more detail and can be more easily

taken into consideration and addressed (e.g., when designing benchmarks). Finally, we believe that this survey will allow future works to determine and compare how the area of addressing the sensitivity to the effects of randomness is continuously advancing.

Acknowledgments

This research was partially supported by TAILOR, a project funded by EU Horizon 2020 research and innovation programme under GA No. [952215](#); DisAI, a project funded by European Union under the Horizon Europe, GA No. [101079164](#); and by vera.ai project funded by the European Union under the Horizon Europe, GA No. [101070093](#).

References

- [1] Rishabh Adiga, Lakshminarayanan Subramanian, and Varun Chandrasekaran. 2024. Designing Informative Metrics for Few-Shot Example Selection. *arXiv preprint arXiv:2403.03861* (2024).
- [2] Mayank Agarwal, Mikhail Yurochkin, and Yuekai Sun. 2021. On sensitivity of meta-learning to support data. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., 20447–20460.
- [3] Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. In-context Examples Selection for Machine Translation. In *Findings of the Association for Computational Linguistics: ACL 2023*. ACL, Toronto, Canada, 8857–8873. <https://doi.org/10.18653/v1/2023.findings-acl.564>
- [4] Anirudh Ajith, Mengzhou Xia, Ameet Deshpande, and Karthik R Narasimhan. 2023. InstructEval: Systematic Evaluation of Instruction Selection Methods. In *Ro-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*. <https://openreview.net/forum?id=6FwaSOEeKD>
- [5] Riccardo Albertoni, Sara Colantonio, Piotr Skrzypczyński, and Jerzy Stefanowski. 2023. Reproducibility of Machine Learning: Terminology, Recommendations and Open Issues. *arXiv preprint arXiv:2302.12691* (2023).
- [6] James Urquhart Allingham, Jie Ren, Michael W Dusenberry, Xiuye Gu, Yin Cui, Dustin Tran, Jeremiah Zhe Liu, and Balaji Lakshminarayanan. 2023. A Simple Zero-shot Prompt Weighting Technique to Improve Prompt Ensembling in Text-Image Models. In *Proceedings of the 40th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 202)*. PMLR, 547–568. <https://proceedings.mlr.press/v202/allingham23a.html>
- [7] Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, et al. 2024. When benchmarks are targets: Revealing the sensitivity of large language model leaderboards. *arXiv preprint arXiv:2402.01781* (2024).
- [8] Shengnan An, Bo Zhou, Zeqi Lin, Qiang Fu, Bei Chen, Nanning Zheng, Weizhu Chen, and Jian-Guang Lou. 2023. Skill-Based Few-Shot Selection for In-Context Learning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. ACL, Singapore, 13472–13492. <https://doi.org/10.18653/v1/2023.emnlp-main.831>
- [9] Alan Ansell, Edoardo Maria Ponti, Anna Korhonen, and Ivan Vulić. 2021. Composable sparse fine-tuning for cross-lingual transfer. *arXiv preprint arXiv:2110.07560* (2021).
- [10] Shubham Atreja, Joshua Ashkinaze, Lingyao Li, Julia Mendelsohn, and Libby Hemphill. 2024. Prompt Design Matters for Computational Social Science Tasks but in Unpredictable Ways. *arXiv preprint arXiv:2406.11980* (2024).
- [11] Sinjini Banerjee, Tim Marrinan, Reilly Cannon, Tony Chiang, and Anand D. Sarwate. 2024. Measuring model variability using robust non-parametric testing. *arXiv:2406.08307 [stat.ML]* <https://arxiv.org/abs/2406.08307>
- [12] Trapit Bansal, Rishikesh Jha, and Andrew McCallum. 2020. Learning to Few-Shot Learn Across Diverse Natural Language Classification Tasks. In *Proceedings of the 28th International Conference on Computational Linguistics*. 5108–5123.
- [13] Amanda Bertsch, Maor Ivgi, Uri Alon, Jonathan Berant, Matthew R Gormley, and Graham Neubig. 2024. In-context learning with long-context models: An in-depth exploration. *arXiv preprint arXiv:2405.00200* (2024).
- [14] Stella Biderman, Hailey Schoelkopf, Lintang Sutawika, Leo Gao, Jonathan Tow, Baber Abbasi, Alham Fikri Aji, Pawan Sasanka Ammanamanchi, Sidney Black, Jordan Clive, et al. 2024. Lessons from the Trenches on Reproducible Evaluation of Language Models. *arXiv preprint arXiv:2405.14782* (2024).
- [15] Thomas Boquet, Laure Delisle, Denis Kochetkov, Nathan Schucher, Boris N Oreshkin, and Julien Cornebise. 2019. Reproducibility and Stability Analysis in Metric-Based Few-Shot Learning. In *RML@ ICLR*, Vol. 3. <https://openreview.net/forum?id=B1g-SnUaUN>
- [16] Xavier Bouthillier, Pierre Delaunay, Mirko Bronzi, Assya Trofimov, Brennan Nichyporuk, Justin Szeto, Naz Sepah, Edward Raff, Kanika Madan, Vikram Voleti, Samira Ebrahimi Kahou, Vincent Michalski, Dmitriy Serdyuk, Tal Arbel, Chris Pal, Gaël Varoquaux, and Pascal Vincent. 2021. Accounting for Variance in Machine Learning Benchmarks. In *Proceedings of Machine Learning and Systems*, Vol. 3. 747–769. <https://proceedings.mlsys.org/paper/2021/hash/cfecdb276f634854f3ef915e2e980c31-Abstract.html>

- [17] Xavier Bouthillier, César Laurent, and Pascal Vincent. 2019. Unreproducible Research is Reproducible. In *Proc. of the 36th International Conf. on Machine Learning*. PMLR, 725–734. <https://proceedings.mlr.press/v97/bouthillier19a.html>
- [18] Jonathan Bragg, Arman Cohan, Kyle Lo, and Iz Beltagy. 2021. FLEX: Unifying Evaluation for Few-Shot NLP. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., 15787–15800.
- [19] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D 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 Ziegler, Jeffrey Wu, Clemens Winter, Chris 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. In *Advances in Neural Information Processing Systems*, Vol. 33. Curran Associates, Inc., 1877–1901.
- [20] Ernie Chang, Xiaoyu Shen, Hui-Syuan Yeh, and Vera Demberg. 2021. On Training Instance Selection for Few-Shot Neural Text Generation. In *Proc. of the 59th Annual Meeting of the ACL and the 11th International Joint Conference on Natural Language Processing*. ACL, Online, 8–13. <https://doi.org/10.18653/v1/2021.acl-short.2>
- [21] Haw-Shiuan Chang, Ruei-Yao Sun, Kathryn Ricci, and Andrew McCallum. 2023. Multi-CLS BERT: An Efficient Alternative to Traditional Ensembling. In *Proc. of the 61st Annual Meeting of the ACL*. ACL, Toronto, Canada, 821–854. <https://doi.org/10.18653/v1/2023.acl-long.48>
- [22] Ting-Yun Chang and Robin Jia. 2023. Data Curation Alone Can Stabilize In-context Learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. ACL, Toronto, Canada, 8123–8144. <https://doi.org/10.18653/v1/2023.acl-long.452>
- [23] Boyuan Chen, Mingzhi Wen, Yong Shi, Dayi Lin, Gopi Krishnan Rajbahadur, and Zhen Ming (Jack) Jiang. 2022. Towards training reproducible deep learning models. In *Proceedings of the 44th International Conference on Software Engineering (Pittsburgh, Pennsylvania) (ICSE '22)*. Association for Computing Machinery, New York, NY, USA, 2202–2214. <https://doi.org/10.1145/3510003.3510163>
- [24] Derek Chen, Kun Qian, and Zhou Yu. 2023. Stabilized In-Context Learning with Pre-trained Language Models for Few Shot Dialogue State Tracking. In *Findings of the Association for Computational Linguistics: EACL 2023*. ACL, Dubrovnik, Croatia, 1551–1564. <https://doi.org/10.18653/v1/2023.findings-eacl.115>
- [25] Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. 2022. Revisiting Parameter-Efficient Tuning: Are We Really There Yet?. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. ACL, Abu Dhabi, United Arab Emirates, 2612–2626. <https://doi.org/10.18653/v1/2022.emnlp-main.168>
- [26] Lichang Chen, Jiuhai Chen, Heng Huang, and Minhao Cheng. 2023. PTP: Boosting Stability and Performance of Prompt Tuning with Perturbation-Based Regularizer. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. ACL, Singapore, 13512–13525. <https://doi.org/10.18653/v1/2023.emnlp-main.833>
- [27] Hyunsoo Cho, Hyuhng Joon Kim, Junyeob Kim, Sang-Woo Lee, Sang-goo Lee, Kang Min Yoo, and Taeuk Kim. 2023. Prompt-augmented linear probing: Scaling beyond the limit of few-shot in-context learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 12709–12718.
- [28] Alexandru Cioba, Michael Bromberg, Qian Wang, Ritwik Niyogi, Georgios Batzolis, Jezabel Garcia, Da-shan Shiu, and Alberto Bernacchia. 2022. How to Distribute Data across Tasks for Meta-Learning?. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 6394–6401. <https://doi.org/10.1609/aaai.v36i6.20590>
- [29] Ashok Cutkosky and Francesco Orabona. 2019. Momentum-Based Variance Reduction in Non-Convex SGD. In *Advances in Neural Information Processing Systems*, Vol. 32. Curran Associates, Inc.
- [30] Alexander D'Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D. Hoffman, Farhad Hormozdizari, Neil Houlsby, Shaobo Hou, Ghassen Jerfel, Alan Karthikesalingam, Mario Lucic, Yian Ma, Cory McLean, Diana Mincu, Akinori Mitani, Andrea Montanari, Zachary Nado, Vivek Natarajan, Christopher Nielson, Thomas F. Osborne, Rajiv Raman, Kim Ramasamy, Rory Sayres, Jessica Schrouff, Martin Seneviratne, Shannon Sequeira, Harini Suresh, Victor Veitch, Max Vladymyrov, Xuezhi Wang, Kellie Webster, Steve Yadlowsky, Taedong Yun, Xiaohua Zhai, and D. Sculley. 2022. Underspecification Presents Challenges for Credibility in Modern Machine Learning. *Journal of Machine Learning Research* 23, 226 (2022), 1–61. <http://jmlr.org/papers/v23/20-1335.html>
- [31] Yann N Dauphin and Samuel Schoenholz. 2019. MetaInit: Initializing learning by learning to initialize. In *Advances in Neural Information Processing Systems*, Vol. 32. Curran Associates, Inc.
- [32] Mostafa Dehghani, Yi Tay, Alexey A. Gritsenko, Zhe Zhao, Neil Houlsby, Fernando Diaz, Donald Metzler, and Oriol Vinyals. 2021. The Benchmark Lottery. <https://doi.org/10.48550/arXiv.2107.07002>
- [33] 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. <https://doi.org/10.48550/arXiv.2002.06305>
- [34] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A Survey for In-context Learning. *arXiv preprint arXiv:2301.00234* (2022).

- [35] Rotem Dror, Segev Shlomov, and Roi Reichart. 2019. Deep Dominance - How to Properly Compare Deep Neural Models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. ACL, Florence, Italy, 2773–2785. <https://doi.org/10.18653/v1/P19-1266>
- [36] Avia Efrat, Or Honovich, and Omer Levy. 2022. LMentry: A Language Model Benchmark of Elementary Language Tasks. <http://arxiv.org/abs/2211.02069>
- [37] Yu Fei, Yifan Hou, Zeming Chen, and Antoine Bosselut. 2023. Mitigating Label Biases for In-context Learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. ACL, Toronto, Canada, 14014–14031. <https://doi.org/10.18653/v1/2023.acl-long.783>
- [38] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*. PMLR, 1126–1135.
- [39] Minghao Fu, Yun-Hao Cao, and Jianxin Wu. 2022. Worst Case Matters for Few-Shot Recognition. In *European Conference on Computer Vision*. Springer, 99–115.
- [40] Zihao Fu, Haoran Yang, Anthony Man-Cho So, Wai Lam, Lidong Bing, and Nigel Collier. 2023. On the effectiveness of parameter-efficient fine-tuning. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 37. 12799–12807.
- [41] Chengguang Gan and Tatsunori Mori. 2023. Sensitivity and Robustness of Large Language Models to Prompt Template in Japanese Text Classification Tasks. In *Proc. of the 37th Pacific Asia Conference on Language, Information and Computation*. ACL, Hong Kong, China, 1–11. <https://aclanthology.org/2023.paclic-1.1>
- [42] Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making Pre-trained Language Models Better Few-shot Learners. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*. 3816–3830.
- [43] Hila Gonen, Srinu Iyer, Terra Blevins, Noah Smith, and Luke Zettlemoyer. 2023. Demystifying Prompts in Language Models via Perplexity Estimation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. ACL, Singapore, 10136–10148. <https://doi.org/10.18653/v1/2023.findings-emnlp.679>
- [44] Odd Erik Gundersen, Kevin Coakley, and Christine Kirkpatrick. 2022. Sources of Irreproducibility in Machine Learning: A Review. *arXiv preprint arXiv:2204.07610* (2022).
- [45] Odd Erik Gundersen, Saeid Shamsaliei, Håkon Sletten Kjærnl, and Helge Langseth. 2023. On Reporting Robust and Trustworthy Conclusions from Model Comparison Studies Involving Neural Networks and Randomness. In *Proceedings of the 2023 ACM Conference on Reproducibility and Replicability (ACM REP '23)*. Association for Computing Machinery, New York, NY, USA, 37–61. <https://doi.org/10.1145/3589806.3600044>
- [46] Vipul Gupta, David Pantoja, Candace Ross, Adina Williams, and Megan Ung. 2024. Changing answer order can decrease mmlu accuracy. *arXiv preprint arXiv:2406.19470* (2024).
- [47] Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. 2024. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608* (2024).
- [48] Christopher Hidey, Fei Liu, and Rahul Goel. 2022. Reducing Model Churn: Stable Re-training of Conversational Agents. In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*. ACL, Edinburgh, UK, 14–25. <https://aclanthology.org/2022.sigdial-1.2>
- [49] Timothy Hospedales, Andreas Antoniou, Paul Micaelli, and Amos Storkey. 2021. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence* 44, 9 (2021), 5149–5169.
- [50] Yutai Hou, Hongyuan Dong, Xinghao Wang, Bohan Li, and Wanxiang Che. 2022. MetaPrompting: Learning to Learn Better Prompts. In *Proc. of the 29th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Gyeongju, Republic of Korea, 3251–3262. <https://aclanthology.org/2022.coling-1.287>
- [51] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *International conference on machine learning*. PMLR, 2790–2799.
- [52] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=nZeVKeeFYf9>
- [53] Hang Hua, Xingjian Li, Dejiong Dou, Chengzhong Xu, and Jiebo Luo. 2021. Noise Stability Regularization for Improving BERT Fine-tuning. In *Proc. of the 2021 Conference of the NAACL: Human Language Technologies*. ACL, Online, 3229–3241. <https://doi.org/10.18653/v1/2021.naacl-main.258>
- [54] Baijun Ji, Xiangyu Duan, Zhenyu Qiu, Tong Zhang, Junhui Li, Hao Yang, and Min Zhang. 2024. Submodular-based In-context Example Selection for LLMs-based Machine Translation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*. ELRA and ICCL, Torino, Italia, 15398–15409. <https://aclanthology.org/2024.lrec-main.1337>
- [55] Mingjian Jiang, Yangjun Ruan, Sicong Huang, Saifei Liao, Silviu Pitis, Roger Baker Grosse, and Jimmy Ba. 2023. Calibrating Language Models via Augmented Prompt Ensembles. (2023).

- [56] Iman Jundi and Gabriella Lapesa. 2022. How to Translate Your Samples and Choose Your Shots? Analyzing Translate-train & Few-shot Cross-lingual Transfer. In *Findings of the Assoc. for Comp. Linguistics: NAACL 2022*. ACL, Seattle, United States, 129–150. <https://doi.org/10.18653/v1/2022.findings-naacl.11>
- [57] Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. 2023. Challenges and applications of large language models. *arXiv preprint arXiv:2307.10169* (2023).
- [58] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL-HLT*. 4171–4186.
- [59] Urja Khurana, Eric Nalisnick, and Antske Fokkens. 2021. How Emotionally Stable is ALBERT? Testing Robustness with Stochastic Weight Averaging on a Sentiment Analysis Task. In *Proc. of the 2nd Workshop on Evaluation and Comparison of NLP Systems*. ACL, Punta Cana, Dominican Republic, 16–31. <https://doi.org/10.18653/v1/2021.eval4nlp-1.3>
- [60] Abdullatif Köksal, Timo Schick, and Hinrich Schuetze. 2023. MEAL: Stable and Active Learning for Few-Shot Prompting. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. ACL, Singapore, 506–517. <https://doi.org/10.18653/v1/2023.findings-emnlp.36>
- [61] Sawan Kumar and Partha Talukdar. 2021. Reordering Examples Helps during Priming-based Few-Shot Learning. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. ACL, Online, 4507–4518. <https://doi.org/10.18653/v1/2021.findings-acl.395>
- [62] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. In *International Conference on Learning Representations*.
- [63] Hung-yi Lee, Shang-Wen Li, and Thang Vu. 2022. Meta Learning for Natural Language Processing: A Survey. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Seattle, United States, 666–684. <https://doi.org/10.18653/v1/2022.naacl-main.49>
- [64] Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and Chelsea Finn. 2022. Surgical Fine-Tuning Improves Adaptation to Distribution Shifts. In *NeurIPS 2022 Workshop on Distribution Shifts: Connecting Methods and Applications*.
- [65] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 3045–3059. <https://doi.org/10.18653/v1/2021.emnlp-main.243>
- [66] Kfir Levy, Ali Kavis, and Volkan Cevher. 2021. STORM+: Fully Adaptive SGD with Recursive Momentum for Nonconvex Optimization. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., 20571–20582.
- [67] Bohan Li, Longxu Dou, Yutai Hou, Yunlong Feng, Honglin Mu, and Wanxiang Che. 2023. MixPro: Simple yet Effective Data Augmentation for Prompt-based Learning. *arXiv preprint arXiv:2304.09402* (2023).
- [68] Chen Li, Yixiao Ge, Dian Li, and Ying Shan. 2024. Vision-Language Instruction Tuning: A Review and Analysis. *Transactions on Machine Learning Research* (2024). <https://openreview.net/forum?id=ul2tbUPtHQ> Survey Certification.
- [69] Jia Li, Ge Li, Chongyang Tao, Huangzhao Zhang, Fang Liu, and Zhi Jin. 2023. Large Language Model-Aware In-Context Learning for Code Generation. *arXiv preprint arXiv:2310.09748* (2023).
- [70] Xiaonan Li and Xipeng Qiu. 2023. Finding Support Examples for In-Context Learning. In *Findings of the ACL: EMNLP 2023*. ACL, Singapore, 6219–6235. <https://doi.org/10.18653/v1/2023.findings-emnlp.411>
- [71] Xiang Lisa Li and Percy Liang. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*. Association for Computational Linguistics, Online, 4582–4597. <https://doi.org/10.18653/v1/2021.acl-long.353>
- [72] Chen Liang, Pengcheng He, Yelong Shen, Weizhu Chen, and Tuo Zhao. 2022. CAMERO: Consistency Regularized Ensemble of Perturbed Language Models with Weight Sharing. In *Proc. of the 60th Annual Meeting of the Association for Computational Linguistics*. ACL, Dublin, Ireland, 7162–7175. <https://doi.org/10.18653/v1/2022.acl-long.495>
- [73] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Alexander Cosgrove, Christopher D Manning, Christopher Re, Diana Acosta-Navas, Drew Arad Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023. Holistic Evaluation of Language Models. *Transactions on Machine Learning Research* (2023). <https://openreview.net/forum?id=iO4LZibEqW> Featured Certification, Expert Certification.

- [74] Hongfu Liu and Ye Wang. 2023. Towards Informative Few-Shot Prompt with Maximum Information Gain for In-Context Learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. ACL, Singapore, 15825–15838. <https://doi.org/10.18653/v1/2023.findings-emnlp.1060>
- [75] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. What Makes Good In-Context Examples for GPT-3?. In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*. ACL, Dublin, Ireland and Online, 100–114. <https://doi.org/10.18653/v1/2022.deelio-1.10>
- [76] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-Train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. *ACM Comput. Surv.* 55, 9, Article 195 (jan 2023), 35 pages. <https://doi.org/10.1145/3560815>
- [77] Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2023. What Makes Good Data for Alignment? A Comprehensive Study of Automatic Data Selection in Instruction Tuning. <https://openreview.net/forum?id=BTKAeLqLMw>
- [78] Xiaoming Liu, Chen Liu, Zhaohan Zhang, Chengzhengxu Li, Longtian Wang, Yu Lan, and Chao Shen. 2024. StablePT: Towards Stable Prompting for Few-shot Learning via Input Separation. *arXiv preprint arXiv:2404.19335* (2024).
- [79] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2023. GPT understands, too. *AI Open* (2023). <https://doi.org/10.1016/j.aiopen.2023.08.012>
- [80] Yinpeng Liu, Jiawei Liu, Xiang Shi, Qikai Cheng, and Wei Lu. 2024. Let’s Learn Step by Step: Enhancing In-Context Learning Ability with Curriculum Learning. *arXiv preprint arXiv:2402.10738* (2024).
- [81] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [82] Robert Logan IV, Ivana Balazevic, Eric Wallace, Fabio Petroni, Sameer Singh, and Sebastian Riedel. 2022. Cutting Down on Prompts and Parameters: Simple Few-Shot Learning with Language Models. In *Findings of the Association for Computational Linguistics: ACL 2022*. Association for Computational Linguistics, Dublin, Ireland, 2824–2835. <https://doi.org/10.18653/v1/2022.findings-acl.222>
- [83] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically Ordered Prompts and Where to Find Them: Overcoming Few-Shot Prompt Order Sensitivity. In *Proc. of the 60th Annual Meeting of the ACL*. ACL, Dublin, Ireland, 8086–8098. <https://doi.org/10.18653/v1/2022.acl-long.556>
- [84] Huan Ma, Changqing Zhang, Yatao Bian, Lemaou Liu, Zhirui Zhang, Peilin Zhao, Shu Zhang, Huazhu Fu, Qinghua Hu, and Bingzhe Wu. 2023. Fairness-guided Few-shot Prompting for Large Language Models. In *Advances in Neural Information Processing Systems*, Vol. 36. Curran Associates, Inc., 43136–43155. https://proceedings.neurips.cc/paper_files/paper/2023/file/8678da90126aa58326b2fc0254b33a8c-Paper-Conference.pdf
- [85] Yubo Ma, Yixin Cao, Yong Hong, and Aixin Sun. 2023. Large Language Model Is Not a Good Few-shot Information Extractor, but a Good Reranker for Hard Samples!. In *Findings of ACL: EMNLP 2023*. ACL, Singapore, 10572–10601. <https://doi.org/10.18653/v1/2023.findings-emnlp.710>
- [86] Lovish Madaan, Aaditya K Singh, Rylan Schaeffer, Andrew Poulton, Sanmi Koyejo, Pontus Stenetorp, Sharan Narang, and Dieuwke Hupkes. 2024. Quantifying Variance in Evaluation Benchmarks. *arXiv preprint arXiv:2406.10229* (2024).
- [87] Katerina Margatina, Timo Schick, Nikolaos Aletras, and Jane Dwivedi-Yu. 2023. Active Learning Principles for In-Context Learning with Large Language Models. In *Findings of ACL: EMNLP 2023*. ACL, Singapore, 5011–5034. <https://doi.org/10.18653/v1/2023.findings-emnlp.334>
- [88] Costas Mavromatis, Balasubramaniam Srinivasan, Zhengyuan Shen, Jiani Zhang, Huzefa Rangwala, Christos Faloutsos, and George Karypis. 2023. Which examples to annotate for in-context learning? towards effective and efficient selection. *arXiv preprint arXiv:2310.20046* (2023).
- [89] R. Thomas McCoy, Junghyun Min, and Tal Linzen. 2020. BERTs of a feather do not generalize together: Large variability in generalization across models with similar test set performance. In *Proc. of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*. ACL, Online, 217–227. <https://doi.org/10.18653/v1/2020.blackboxnlp-1.21>
- [90] Yu Meng, Martin Michalski, Jiaxin Huang, Yu Zhang, Tarek Abdelzaher, and Jiawei Han. 2022. Tuning Language Models as Training Data Generators for Augmentation-Enhanced Few-Shot Learning. <http://arxiv.org/abs/2211.03044>
- [91] Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Noisy Channel Language Model Prompting for Few-Shot Text Classification. In *Proc. of the 60th Annual Meeting of the Association for Computational Linguistics*. ACL, Dublin, Ireland, 5316–5330. <https://doi.org/10.18653/v1/2022.acl-long.365>
- [92] Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. 2023. State of what art? a call for multi-prompt llm evaluation. *arXiv preprint arXiv:2401.00595* (2023).
- [93] David Moher, Alessandro Liberati, Jennifer Tetzlaff, Douglas G Altman, and PRISMA Group*. 2009. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *Annals of internal medicine* 151, 4 (2009),

- 264–269.
- [94] Marius Mosbach, Maksym Andriushchenko, and Dietrich Klakow. 2021. On the Stability of Fine-tuning BERT: Misconceptions, Explanations, and Strong Baselines. <https://openreview.net/forum?id=nzpl.WnVAyah>
 - [95] Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. 2023. Few-shot Fine-tuning vs. In-context Learning: A Fair Comparison and Evaluation. In *Findings of the ACL: ACL 2023*. ACL, Toronto, Canada, 12284–12314. <https://doi.org/10.18653/v1/2023.findings-acl.779>
 - [96] Subhabrata Mukherjee, Xiaodong Liu, Guoqing Zheng, Saghar Hosseini, Hao Cheng, Greg Yang, Christopher Meek, Ahmed Hassan Awadallah, and Jianfeng Gao. 2021. CLUES: Few-Shot Learning Evaluation in Natural Language Understanding. <http://arxiv.org/abs/2111.02570>
 - [97] Tai Nguyen and Eric Wong. 2023. In-context Example Selection with Influences. *arXiv preprint arXiv:2302.11042* (2023).
 - [98] Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999* (2018).
 - [99] Uche Osahor and Nasser M. Nasrabadi. 2022. Ortho-Shot: Low Displacement Rank Regularization with Data Augmentation for Few-Shot Learning. In *2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*. 2040–2049. <https://doi.org/10.1109/WACV51458.2022.00210>
 - [100] Matthew J Page, Joanne E McKenzie, Patrick M Bossuyt, Isabelle Boutron, Tammy C Hoffmann, Cynthia D Mulrow, Larissa Shamseer, Jennifer M Tetzlaff, Elie A Akl, Sue E Brennan, et al. 2021. The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. *Systematic reviews* 10, 1 (2021), 1–11.
 - [101] Kaihang Pan, Juncheng Li, Hongye Song, Jun Lin, Xiaozhong Liu, and Siliang Tang. 2023. Self-supervised Meta-Prompt Learning with Meta-Gradient Regularization for Few-shot Generalization. In *Findings of the ACL: EMNLP 2023*. ACL, Singapore, 1059–1077. <https://doi.org/10.18653/v1/2023.findings-emnlp.75>
 - [102] Branislav Pecher, Jan Cegin, Robert Belanec, Jakub Simko, Ivan Srba, and Maria Bielikova. 2024. Fighting Randomness with Randomness: Mitigating Optimisation Instability of Fine-Tuning using Delayed Ensemble and Noisy Interpolation. *arXiv preprint arXiv:2406.12471* (2024).
 - [103] Branislav Pecher, Ivan Srba, and Maria Bielikova. 2024. Comparing Specialised Small and General Large Language Models on Text Classification: 100 Labelled Samples to Achieve Break-Even Performance. *arXiv preprint arXiv:2402.12819* (2024).
 - [104] Branislav Pecher, Ivan Srba, and Maria Bielikova. 2024. On Sensitivity of Learning with Limited Labelled Data to the Effects of Randomness: Impact of Interactions and Systematic Choices. *arXiv preprint arXiv:2402.12817* (2024).
 - [105] Branislav Pecher, Ivan Srba, Maria Bielikova, and Joaquin Vanschoren. 2024. Automatic Combination of Sample Selection Strategies for Few-Shot Learning. *arXiv preprint arXiv:2402.03038* (2024).
 - [106] Keqin Peng, Liang Ding, Yancheng Yuan, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2024. Revisiting Demonstration Selection Strategies in In-Context Learning. *arXiv preprint arXiv:2401.12087* (2024).
 - [107] Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. True Few-Shot Learning with Language Models. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., 11054–11070.
 - [108] Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of options in multiple-choice questions. *arXiv preprint arXiv:2308.11483* (2023).
 - [109] Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. AdapterFusion: Non-Destructive Task Composition for Transfer Learning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics, Online, 487–503. <https://doi.org/10.18653/v1/2021.eacl-main.39>
 - [110] Hung Viet Pham, Shangshu Qian, Jiannan Wang, Thibaud Lutellier, Jonathan Rosenthal, Lin Tan, Yaoliang Yu, and Nachiappan Nagappan. 2021. Problems and opportunities in training deep learning software systems: an analysis of variance. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering (ASE '20)*. Association for Computing Machinery, New York, NY, USA, 771–783. <https://doi.org/10.1145/3324884.3416545>
 - [111] Jason Phang, Thibault Févry, and Samuel R. Bowman. 2019. Sentence Encoders on STILTs: Supplementary Training on Intermediate Labeled-data Tasks. <https://doi.org/10.48550/arXiv.1811.01088>
 - [112] Felipe Maia Polo, Ronald Xu, Lucas Weber, Mirian Silva, Onkar Bhardwaj, Leshem Choshen, Allysson Flavio Melo de Oliveira, Yuekai Sun, and Mikhail Yurochkin. 2024. Efficient multi-prompt evaluation of LLMs. *arXiv preprint arXiv:2405.17202* (2024).
 - [113] Reid Pryzant, Dan Iter, Jerry Li, Yin Lee, Chenguang Zhu, and Michael Zeng. 2023. Automatic Prompt Optimization with “Gradient Descent” and Beam Search. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. ACL, Singapore, 7957–7968. <https://doi.org/10.18653/v1/2023.emnlp-main.494>
 - [114] Jian Qian, Miao Sun, Sifan Zhou, Ziyu Zhao, Ruizhi Hun, and Patrick Chiang. 2024. Sub-SA: Strengthen In-context Learning via Submodular Selective Annotation. *arXiv:2407.05693* [cs.LG] <https://arxiv.org/abs/2407.05693>

- [115] Chengwei Qin, Shafiq Joty, Qian Li, and Ruochen Zhao. 2023. Learning to Initialize: Can Meta Learning Improve Cross-task Generalization in Prompt Tuning?. In *Proc. of the 61st Annual Meeting of the ACL*. ACL, Toronto, Canada, 11802–11832. <https://doi.org/10.18653/v1/2023.acl-long.659>
- [116] Chengwei Qin, Aston Zhang, Anirudh Dagar, and Wenming Ye. 2023. In-context learning with iterative demonstration selection. *arXiv preprint arXiv:2310.09881* (2023).
- [117] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. 2019. Rapid Learning or Feature Reuse? Towards Understanding the Effectiveness of MAML. In *International Conference on Learning Representations*.
- [118] Anastasiia Razdaibiedina, Yuning Mao, Madian Khabsa, Mike Lewis, Rui Hou, Jimmy Ba, and Amjad Almahairi. 2023. Residual Prompt Tuning: improving prompt tuning with residual reparameterization. In *Findings of the ACL: ACL 2023*. ACL, Toronto, Canada, 6740–6757. <https://doi.org/10.18653/v1/2023.findings-acl.421>
- [119] Yuval Reif and Roy Schwartz. 2024. Beyond Performance: Quantifying and Mitigating Label Bias in LLMs. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. ACL, Mexico City, Mexico, 6784–6798. <https://aclanthology.org/2024.naacl-long.378>
- [120] Nils Reimers and Iryna Gurevych. 2017. Reporting Score Distributions Makes a Difference: Performance Study of LSTM-networks for Sequence Tagging. In *Proc. of the 2017 Conference on Empirical Methods in Natural Language Processing*. ACL, Copenhagen, Denmark, 338–348. <https://doi.org/10.18653/v1/D17-1035>
- [121] Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in bertology: What we know about how bert works. *Transactions of the Association for Computational Linguistics* 8 (2020), 842–866.
- [122] Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning To Retrieve Prompts for In-Context Learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. ACL, Seattle, United States, 2655–2671. <https://doi.org/10.18653/v1/2022.naacl-main.191>
- [123] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).
- [124] Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesh Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. 2022. Multitask Prompted Training Enables Zero-Shot Task Generalization. In *International Conf. on Learning Representations*. <https://openreview.net/forum?id=9Vrb9D0Wl4>
- [125] Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose Opinions Do Language Models Reflect?. In *Proceedings of the 40th International Conference on Machine Learning*. PMLR, 29971–30004. <https://proceedings.mlr.press/v202/santurkar23a.html> ISSN: 2640-3498.
- [126] Timo Schick and Hinrich Schütze. 2021. It’s Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. ACL, Online, 2339–2352. <https://doi.org/10.18653/v1/2021.naacl-main.185>
- [127] Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying Language Models’ Sensitivity to Spurious Features in Prompt Design or: How I learned to start worrying about prompt formatting. <https://openreview.net/forum?id=RIu5lyNXJT>
- [128] Thibault Sellam, Steve Yadlowsky, Ian Tenney, Jason Wei, Naomi Saphra, Alexander D’Amour, Tal Linzen, Jasmijn Bastings, Iulia Turc, Jacob Eisenstein, Dipanjan Das, and Ellie Pavlick. 2022. The MultiBERTs: BERT Reproductions for Robustness Analysis. 30. https://openreview.net/forum?id=K0E_F0gFDgA
- [129] Amrith Setlur, Oscar Li, and Virginia Smith. 2021. Is Support Set Diversity Necessary for Meta-Learning? <http://arxiv.org/abs/2011.14048>
- [130] Amrith Setlur, Oscar Li, and Virginia Smith. 2021. Two Sides of Meta-Learning Evaluation: In vs. Out of Distribution. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., 3770–3783.
- [131] Kashun Shum, Shizhe Diao, and Tong Zhang. 2023. Automatic Prompt Augmentation and Selection with Chain-of-Thought from Labeled Data. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. ACL, Singapore, 12113–12139. <https://doi.org/10.18653/v1/2023.findings-emnlp.811>
- [132] Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical Networks for Few-shot Learning. In *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc.
- [133] Sai Ashish Somayajula, Youwei Liang, Li Zhang, Abhishek Singh, and Pengtao Xie. 2024. Generalizable and Stable Finetuning of Pretrained Language Models on Low-Resource Texts. In *Proceedings of the 2024 Conference of the NAACL: Human Language Technologies*. ACL, Mexico City, Mexico, 4936–4953. <https://aclanthology.org/2024.naacl-long.277>
- [134] Yisheng Song, Ting Wang, Puyu Cai, Subrota K Mondal, and Jyoti Prakash Sahoo. 2023. A Comprehensive Survey of Few-Shot Learning: Evolution, Applications, Challenges, and Opportunities. *ACM Comput. Surv.* (feb 2023). <https://doi.org/10.1145/3582688> Just Accepted.

- [135] Taylor Sorensen, Joshua Robinson, Christopher Rytting, Alexander Shaw, Kyle Rogers, Alexia Delorey, Mahmoud Khalil, Nancy Fulda, and David Wingate. 2022. An Information-theoretic Approach to Prompt Engineering Without Ground Truth Labels. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*. ACL, Dublin, Ireland, 819–862. <https://doi.org/10.18653/v1/2022.acl-long.60>
- [136] Michal Štefánik, Marek Kadlčík, Piotr Gramacki, and Petr Sojka. 2023. Resources and Few-shot Learners for In-context Learning in Slavic Languages. In *Proc. of the 9th Workshop on Slavic Natural Language Processing 2023 (SlavicNLP 2023)*. ACL, Dubrovnik, Croatia, 94–105. <https://doi.org/10.18653/v1/2023.bsnlp-1.12>
- [137] Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. 2022. Selective Annotation Makes Language Models Better Few-Shot Learners. <https://openreview.net/forum?id=qY1hlv7gwg>
- [138] Cecilia Summers and Michael J. Dinneen. 2021. Nondeterminism and Instability in Neural Network Optimization. In *Proceedings of the 38th International Conference on Machine Learning*. PMLR, 9913–9922. <https://proceedings.mlr.press/v139/summers21a.html> ISSN: 2640-3498.
- [139] Jiuding Sun, Chantal Shaib, and Byron C. Wallace. 2023. Evaluating the Zero-shot Robustness of Instruction-tuned Language Models. <https://openreview.net/forum?id=g9diuvxN6D>
- [140] Yingjie Tian, Xiaoxi Zhao, and Wei Huang. 2022. Meta-learning approaches for learning-to-learn in deep learning: A survey. *Neurocomputing* 494 (2022), 203–223. <https://doi.org/10.1016/j.neucom.2022.04.078>
- [141] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. 2016. Matching Networks for One Shot Learning. In *Advances in Neural Information Processing Systems*, Vol. 29. Curran Associates, Inc.
- [142] Anton Voronov, Lena Wolf, and Max Ryabinin. 2024. Mind your format: Towards consistent evaluation of in-context learning improvements. *arXiv preprint arXiv:2401.06766* (2024).
- [143] Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou, and Daniel Cer. 2022. SPoT: Better Frozen Model Adaptation through Soft Prompt Transfer. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*. ACL, Dublin, Ireland, 5039–5059. <https://doi.org/10.18653/v1/2022.acl-long.346>
- [144] Risto Vuorio, Shao-Hua Sun, Hexiang Hu, and Joseph J Lim. 2019. Multimodal Model-Agnostic Meta-Learning via Task-Aware Modulation. In *Advances in Neural Information Processing Systems*, Vol. 32. Curran Associates, Inc.
- [145] Lingxiao Wang, Kevin Huang, Tengyu Ma, Quanquan Gu, and Jing Huang. 2021. Variance-reduced First-order Meta-learning for Natural Language Processing Tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. ACL, Online, 2609–2615. <https://doi.org/10.18653/v1/2021.naacl-main.206>
- [146] Lijing Wang, Yingya Li, Timothy Miller, Steven Bethard, and Guergana Savova. 2023. Two-Stage Fine-Tuning for Improved Bias and Variance for Large Pretrained Language Models. In *Proc. of the 61st Annual Meeting of the ACL*. ACL, Toronto, Canada, 15746–15761. <https://doi.org/10.18653/v1/2023.acl-long.877>
- [147] Lucas Weber, Elia Bruni, and Dieuwke Hupkes. 2023. Mind the instructions: a holistic evaluation of consistency and interactions in prompt-based learning. In *Proc. of the 27th Conference on Computational Natural Language Learning (CoNLL)*. ACL, Singapore, 294–313. <https://doi.org/10.18653/v1/2023.conll-1.20>
- [148] Albert Webson and Ellie Pavlick. 2022. Do Prompt-Based Models Really Understand the Meaning of Their Prompts?. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. ACL, Seattle, United States, 2300–2344. <https://doi.org/10.18653/v1/2022.naacl-main.167>
- [149] Sheng-Lun Wei, Cheng-Kuang Wu, Hen-Hsen Huang, and Hsin-Hsi Chen. 2024. Unveiling Selection Biases: Exploring Order and Token Sensitivity in Large Language Models. *arXiv preprint arXiv:2406.03009* (2024).
- [150] Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. 2022. NoisyTune: A Little Noise Can Help You Finetune Pretrained Language Models Better. In *Proc. of the 60th Annual Meeting of the Association for Computational Linguistics*. ACL, Dublin, Ireland, 680–685. <https://doi.org/10.18653/v1/2022.acl-short.76>
- [151] Sherry Wu, Hua Shen, Daniel S Weld, Jeffrey Heer, and Marco Tulio Ribeiro. 2023. ScatterShot: Interactive In-context Example Curation for Text Transformation. In *Proceedings of the 28th International Conference on Intelligent User Interfaces* (Sydney, NSW, Australia) (IUI '23). Association for Computing Machinery, New York, NY, USA, 353–367. <https://doi.org/10.1145/3581641.3584059>
- [152] Shijie Wu, Benjamin Van Durme, and Mark Dredze. 2022. Zero-shot Cross-lingual Transfer is Under-specified Optimization. In *Proceedings of the 7th Workshop on Representation Learning for NLP*. ACL, Dublin, Ireland, 236–248. <https://doi.org/10.18653/v1/2022.repl4nlp-1.25>
- [153] Zhaoxuan Wu, Xiaoqiang Lin, Zhongxiang Dai, Wenyang Hu, Yao Shu, See-Kiong Ng, Patrick Jaillet, and Bryan Kian Hsiang Low. 2024. Prompt Optimization with EASE? Efficient Ordering-aware Automated Selection of Exemplars. In *ICML 2024 Workshop on In-Context Learning*. <https://openreview.net/forum?id=TYxOXHYU6b>
- [154] Zhiyong Wu, Yaoxiang Wang, Jiacheng Ye, and Lingpeng Kong. 2023. Self-Adaptive In-Context Learning: An Information Compression Perspective for In-Context Example Selection and Ordering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. ACL, Toronto, Canada, 1423–1436. <https://doi.org/>

10.18653/v1/2023.acl-long.79

- [155] Patrick Xia, Shijie Wu, and Benjamin Van Durme. 2020. Which *BERT? A Survey Organizing Contextualized Encoders. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 7516–7533. <https://doi.org/10.18653/v1/2020.emnlp-main.608>
- [156] Runxin Xu, Fuli Luo, Zhiyuan Zhang, Chuanqi Tan, Baobao Chang, Songfang Huang, and Fei Huang. 2021. Raise a Child in Large Language Model: Towards Effective and Generalizable Fine-tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. ACL, Online and Punta Cana, Dominican Republic, 9514–9528. <https://doi.org/10.18653/v1/2021.emnlp-main.749>
- [157] Shangqing Xu and Chao Zhang. 2024. Misconfidence-based demonstration selection for llm in-context learning. *arXiv preprint arXiv:2401.06301* (2024).
- [158] Xin Xu, Yue Liu, Panupong Pasupat, Mehran Kazemi, et al. 2024. In-context learning with retrieved demonstrations for language models: A survey. *arXiv preprint arXiv:2401.11624* (2024).
- [159] Zhiyang Xu, Ying Shen, and Lifu Huang. 2023. MultiInstruct: Improving Multi-Modal Zero-Shot Learning via Instruction Tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. ACL, Toronto, Canada, 11445–11465. <https://doi.org/10.18653/v1/2023.acl-long.641>
- [160] Hansi Yang and James Kwok. 2022. Efficient Variance Reduction for Meta-learning. In *Proceedings of the 39th International Conference on Machine Learning*. PMLR, 25070–25095. <https://proceedings.mlr.press/v162/yang22g.html>
- [161] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems* 32 (2019).
- [162] Zhao Yang, Yuanzhe Zhang, Dianbo Sui, Cao Liu, Jun Zhao, and Kang Liu. 2023. Representative Demonstration Selection for In-Context Learning with Two-Stage Determinantal Point Process. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. ACL, Singapore, 5443–5456. <https://doi.org/10.18653/v1/2023.emnlp-main.331>
- [163] Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023. Compositional exemplars for in-context learning. In *International Conference on Machine Learning*. PMLR, 39818–39833.
- [164] Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. 2021. CrossFit: A Few-shot Learning Challenge for Cross-task Generalization in NLP. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. ACL, Online and Punta Cana, Dominican Republic, 7163–7189. <https://doi.org/10.18653/v1/2021.emnlp-main.572>
- [165] Yue Yu, Rongzhi Zhang, Ran Xu, Jieyu Zhang, Jiaming Shen, and Chao Zhang. 2023. Cold-Start Data Selection for Better Few-shot Language Model Fine-tuning: A Prompt-based Uncertainty Propagation Approach. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. ACL, Toronto, Canada, 2499–2521. <https://doi.org/10.18653/v1/2023.acl-long.141>
- [166] Pengwei Zhan, Zhen Xu, Qian Tan, Jie Song, and Ru Xie. 2024. Unveiling the Lexical Sensitivity of LLMs: Combinatorial Optimization for Prompt Enhancement. *arXiv preprint arXiv:2405.20701* (2024).
- [167] Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: a case study. In *Proceedings of the 40th International Conference on Machine Learning (Honolulu, Hawaii, USA) (ICML'23)*. JMLR.org, Article 1722, 19 pages.
- [168] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2018. mixup: Beyond Empirical Risk Minimization. In *International Conference on Learning Representations*.
- [169] Haojie Zhang, Ge Li, Jia Li, Zhongjin Zhang, Yuqi Zhu, and Zhi Jin. 2022. Fine-Tuning Pre-Trained Language Models Effectively by Optimizing Subnetworks Adaptively. 15. <https://openreview.net/forum?id=r6-WNKfyhW>
- [170] Kaiyi Zhang, Ang Lv, Yuhang Chen, Hansen Ha, Tao Xu, and Rui Yan. 2024. Batch-ICL: Effective, Efficient, and Order-Agnostic In-Context Learning. *arXiv preprint arXiv:2401.06469* (2024).
- [171] Miaoran Zhang, Vagrant Gautam, Mingyang Wang, Jesujoba O Alabi, Xiaoyu Shen, Dietrich Klakow, and Marius Mosbach. 2024. The Impact of Demonstrations on Multilingual In-Context Learning: A Multidimensional Analysis. *arXiv preprint arXiv:2402.12976* (2024).
- [172] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792* (2023).
- [173] Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q. Weinberger, and Yoav Artzi. 2021. Revisiting Few-sample BERT Fine-tuning. *arXiv*. <https://openreview.net/forum?id=c01IH43yUF>
- [174] Yiming Zhang, Shi Feng, and Chenhao Tan. 2022. Active Example Selection for In-Context Learning. <http://arxiv.org/abs/2211.04486>
- [175] Zhihan Zhang, Shuohang Wang, Wenhao Yu, Yichong Xu, Dan Iter, Qingkai Zeng, Yang Liu, Chenguang Zhu, and Meng Jiang. 2023. Auto-Instruct: Automatic Instruction Generation and Ranking for Black-Box Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. ACL, Singapore, 9850–9867. <https://doi.org/10.18653/v1/2023.findings-emnlp.659>

- [176] Feng Zhao, Wan Xianlin, Cheng Yan, and Chu Kiong Loo. 2024. Correcting Language Model Bias for Text Classification in True Zero-Shot Learning. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*. ELRA and ICCL, Torino, Italia, 4036–4046. <https://aclanthology.org/2024.lrec-main.359>
- [177] Mengjie Zhao, Yi Zhu, Ehsan Shareghi, Ivan Vulić, Roi Reichart, Anna Korhonen, and Hinrich Schütze. 2021. A Closer Look at Few-Shot Crosslingual Transfer: The Choice of Shots Matters. In *Proc. of the 59th Annual Meeting of the ACL and the 11th International Joint Conference on Natural Language Processing*. ACL, Online, 5751–5767. <https://doi.org/10.18653/v1/2021.acl-long.447>
- [178] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223* (2023).
- [179] Yufeng Zhao, Yoshihiro Sakai, and Naoya Inoue. 2024. NoisyICL: A Little Noise in Model Parameters Calibrates In-context Learning. *arXiv preprint arXiv:2402.05515* (2024).
- [180] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate Before Use: Improving Few-shot Performance of Language Models. In *Proceedings of the 38th International Conference on Machine Learning*. PMLR, 12697–12706. <https://proceedings.mlr.press/v139/zhao21c.html>
- [181] Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. 2023. Large Language Models Are Not Robust Multiple Choice Selectors. <https://openreview.net/forum?id=shr9PXz7T0>
- [182] Yanan Zheng, Jing Zhou, Yujie Qian, Ming Ding, Chonghua Liao, Li Jian, Ruslan Salakhutdinov, Jie Tang, Sebastian Ruder, and Zhilin Yang. 2022. FewNLU: Benchmarking State-of-the-Art Methods for Few-Shot Natural Language Understanding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*. ACL, Dublin, Ireland, 501–516. <https://doi.org/10.18653/v1/2022.acl-long.38>
- [183] Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. 2023. Can chatgpt understand too? a comparative study on chatgpt and fine-tuned bert. *arXiv preprint arXiv:2302.10198* (2023).
- [184] Ruiqi Zhong, Dhruva Ghosh, Dan Klein, and Jacob Steinhardt. 2021. Are Larger Pretrained Language Models Uniformly Better? Comparing Performance at the Instance Level. In *Findings of the Assoc. for Comp. Linguistics: ACL-IJCNLP 2021*. ACL, Online, 3813–3827. <https://doi.org/10.18653/v1/2021.findings-acl.334>
- [185] Chunting Zhou, Junxian He, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2022. Prompt Consistency for Zero-Shot Task Generalization. In *Findings of the ACL: EMNLP 2022*. ACL, Abu Dhabi, United Arab Emirates, 2613–2626. <https://doi.org/10.18653/v1/2022.findings-emnlp.192>
- [186] Han Zhou, Xingchen Wan, Lev Proleev, Diana Mincu, Jilin Chen, Katherine A. Heller, and Subhrajit Roy. 2023. Batch Calibration: Rethinking Calibration for In-Context Learning and Prompt Engineering. <https://openreview.net/forum?id=L3FHM0KZcS>
- [187] Xiang Zhou, Yixin Nie, Hao Tan, and Mohit Bansal. 2020. The Curse of Performance Instability in Analysis Datasets: Consequences, Source, and Suggestions. In *Proc. of the 2020 Conf. on Empirical Methods in Natural Language Processing (EMNLP)*. ACL, Online, 8215–8228. <https://doi.org/10.18653/v1/2020.emnlp-main.659>
- [188] Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Neil Zhenqiang Gong, Yue Zhang, et al. 2023. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts. *arXiv preprint arXiv:2306.04528* (2023).
- [189] Shaolin Zhu, Menglong Cui, and Deyi Xiong. 2024. Towards Robust In-Context Learning for Machine Translation with Large Language Models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*. ELRA and ICCL, Torino, Italia, 16619–16629. <https://aclanthology.org/2024.lrec-main.1444>
- [190] Donglin Zhuang, Xingyao Zhang, Shuaiwen Song, and Sara Hooker. 2022. Randomness in Neural Network Training: Characterizing the Impact of Tooling. *Proceedings of Machine Learning and Systems* 4 (April 2022), 316–336. https://proceedings.mlsys.org/paper_files/paper/2022/hash/427e0e886ebf87538afdf0badb805b7f-Abstract.html
- [191] Yongshuo Zong, Tingyang Yu, Bingchen Zhao, Ruchika Chavhan, and Timothy Hospedales. 2023. Fool your (vision and) language model with embarrassingly simple permutations. *arXiv preprint arXiv:2310.01651* (2023).

A Detailed Paper Categorisation Using the Defined Taxonomy

As a part of the main content of the survey, only a part of the full paper categorisation of the *core* papers according to the defined taxonomy is included, mainly due to its size. The full paper categorisation, along with additional relevant information, such as a basic categorisation of the *recognise* papers is a part of the digital appendix of this survey².

The following information is part of the digital appendix:

- (1) **Full categorisation of *core* papers** – The categorisation is extended with additional taxonomy dimension as well as other metadata. The *machine learning approaches* are further divided into more detail, e.g., in meta-learning, we identify which papers deal with optimisation-based meta-learning and which consider metric-based meta-learning. In addition, we also divide the *tasks performed* to the full list presented in the survey to provide a better overview of which tasks the papers focus on when addressing the effects of randomness. We also add the *dataset* secondary dimension that indicates what datasets the individual papers work with. Besides the changes to the taxonomy dimensions, we also include additional metadata, such as Document Object Identifier (DOI) number (where applicable), year of publication, or link to the paper. Finally, for majority of the papers, we also provide a short description of the main focus for the paper and its findings.
- (2) **Basic categorisation of *recognise* papers** – We include the full list of the *recognise* papers in this digital appendix, along with their metadata (authors, DOI, year of publication). In addition, we identify some of the *randomness factors* the *recognise* papers deal with. Finally, we specify a high-level idea for each paper, for example, whether the paper recognises the problem of instability, and tries to improve the situation.
- (3) **Evolution in number of papers through years** – Using the year of publication for each identified paper, we show the evolution in the number of papers that focus on addressing the effects of randomness.
- (4) **Relation mapping between *randomness factors* and *machine learning approaches* for different *tasks performed*** – As a part of the survey, such mapping was presented in an aggregated form across all the tasks for addressing the effects of randomness. Besides the aggregated mapping, we also include the same table, but for each task separately. This mapping serves as an additional overview of what the individual papers focus on (*randomness factors* in specific *machine learning approaches* and *tasks performed*).

B Description of Machine Learning Approaches

In this part, we provide a basic idea and high-level description of the different *machine learning approaches* that specify our scope in the survey, along with some representatives for each. This should serve as a basic overview for understating the concepts of the approaches required for a better understanding of the surveyed papers, not as a comprehensive overview of the selected approaches. We provide description for *meta-learning*, *language model fine-tuning*, *prompting/in-context learning*, *prompt-based learning* and *parameter-efficient fine-tuning* approaches.

B.1 Meta-Learning

The idea behind *meta-learning* is to learn a general learning algorithm that can generalise across different tasks, with only a few samples available for each task. In essence, meta-learning learns a function (also called *meta-learner*) that takes a set of training data as input and outputs the learned model parameters (also called *base-learner*) that can be used for the specific task. The *few-shot learning* setup is often used, where the samples are distributed across multiple tasks.

²Available at <https://kinit.sk/public/acm-csur-sensitivity-survey.html>

During training, a subset of these tasks is sampled and split into training and testing set. Using the training tasks, the parameters of the *base-learner* are optimised. The parameters of *meta-learner* are then updated based on how good the *base-learner* performs on tasks in the testing set. This training is performed over multiple iterations, while at the start of each iteration, the *base-learner* is initialised or updated based on the information or knowledge contained in the *meta-learner*. At test time, the *base-learner* is primed for the specific task using a few adaptation (or support) samples.

Different approaches of meta-learning differ based on the objective they optimise for and how the *base-learner* and *meta-learner* are defined. In *optimisation-based* meta-learning, the objective is to find an optimal set of parameters in the *meta-learner* that can be quickly adapted to any task when fine-tuned as a part of *base-learner* on a given task. In optimisation learning, the *base-learner* and *meta-learner* often represent the same model, and the *base-learner* is initialised by simply copying all the parameters of *meta-learner*. The optimisation of the *meta-learner* is performed similarly to the supervised learning (where the error is back-propagated through the whole network), but using a second-order derivative. Majority of optimisation-based approaches are based on Model-Agnostic Meta-Learning (MAML) [38], as they are designed to overcome some of the problems present in MAML, for example reducing its computational complexity by approximating the outer loop (e.g., Reptile [98] or First Order MAML [38]) or optimising only part of the network while having the rest pre-trained (e.g. ANIL [117]), or better dealing with task dissimilarity and distribution (e.g. MMAML [144], LEOPARD [12]).

In *metric-based* meta-learning, the objective is to learn good representation for comparing adaptation samples with samples without known labels. The final prediction of the *base-learner* is determined by the class of the most similar adaptation (or support) sample. The most popular approaches for metric-based meta-learning are Prototypical Networks [132], or Matching Networks [141], which only differ in how they function (e.g., comparing based on aggregated information from multiple samples or comparing between each sample separately).

For a more comprehensive overview of the taxonomy of different *meta-learning* approaches and a more detailed description, please refer to [49, 63, 140].

B.2 Language Model Fine-Tuning

The main idea of *language model fine-tuning* is to transfer knowledge from a pre-trained model to a specific task by updating its learned parameters using a few samples. The language model is first trained on large corpora of unrelated texts, or an already pre-trained language model is used. The last layer of such language model is then replaced, or a new one is added for the specific task. Afterwards, the language model is fine-tuned on the few samples available by simply running a training process for a few episodes on the samples with a much lower learning rate.

As the process is pretty straightforward, the approaches differ only in what language model is used (which also defines what corpora it was pretrained on) and what part of the language model is fine-tuned. The common approaches include training all the layers, or only the last few layers of the language model (such as only the added classification layer), or even approaches that fine-tune specific, often disconnected layers [64]. The most popular language models that are used for fine-tuning are BERT [58] and its modifications like RoBERTa [81], ALBERT [62], DistilBERT [123], XLNet [161] and many others.

For a more comprehensive overview of different *language model fine-tuning* approaches, please refer to [121, 155],

B.3 Prompting and In-Context Learning

Prompting (and its related technique in-context learning, also called few-shot prompting) is an emerging paradigm, where a large pre-trained language model is used for different tasks without

first updating its parameters. Instead, all tasks are reformulated as a sequence generation problem and the language model is “prompted” to output a sequence of words. The resulting sequence of words is then mapped to a possible list of labels (e.g., the word “good” is mapped to positive and “terrible” to negative label in sentiment classification). To prime the model for a specific task, the presented input is constructed as a concatenation of instruction for the task, optionally a few labelled samples (serving as context) if in-context learning is used, and a single test sample, for which the label is predicted.

As the final prompt is a concatenation of multiple labelled samples and a single unlabelled one, usually only language models that allow large input sizes can be used. Another important choice in prompting and in-context learning is how the input prompt is designed, either manually, semi-automatically or automatically. The approaches differ on what language models are used, how the prompt format is designed and whether the prompting and in-context learning is combined with some kind of fine-tuning. The most popular model used in in-context learning is GPT-3, with the approaches starting to design automatically generated prompts, although also smaller models are used when combined with fine-tuning [19, 42, 82, 107].

For a more comprehensive overview of different *in-context learning* approaches, please refer to [34, 76, 134].

B.4 Prompt-Based Learning

Prompt-based learning, often also called instruction-tuning, can be viewed as a modification of typical fine-tuning for the large generative language models and their use through prompting and in-context. The idea is to bridge the gap between the next-word prediction objective of the large language models and the objective of the model adhering to human instructions. As such, the goal is to optimise the parameters of the large language models to better follow the instructions that are included as part of the prompts.

Following the typical fine-tuning technique, a pre-trained large language model is further trained on a dataset consisting of pairs of instructions and the outputs, where the instructions denote the instruction in the input prompt and the outputs represent the generated words that are mapped to classes. Such tuning allows for more controllable and predictable model behaviour, improving the capability of the models to perform the specific task by improving the mapping from instructions to the generated words that are further mapped to classes. In addition, it is characterised by the benefits of fine-tuning, i.e., the rapid adaptation to a specific domain and task without extensive retraining or architectural changes.

Even though it provides its benefits, it also faces challenges, especially when it comes to the sensitivity to the effects of randomness. As we are combining the fine-tuning and prompting/in-context learning techniques, it also combines the sensitivity of these models to the randomness factors in fine-tuning and in-context learning. For example, designing high-quality prompts or instructions, the choice of samples, their order, but also the whole optimisation process.

For a more comprehensive overview of different *prompt-based learning* approaches, please refer to [68, 172].

B.5 Parameter-Efficient Fine-Tuning

Parameter-Efficient Fine-Tuning (PEFT) is an extension of regular fine-tuning, where the optimisation of the parameters in the pre-trained large models is done while minimising the number of additional parameters introduced or the computational resources we require. This is particularly important for the massive generative language models with high parameter counts, such as GPT-3, that are often used through prompting or in-context learning. Although PEFT methods are more

popular for such models, they can be utilised for any kind of model to reduce the required resources to run them.

Overall, there are a few categories of PEFT methods, based on how they reduce the number of trainable parameters. First, additive methods, keep all the parameters of the pre-trained models fixed and only add some additional parameters that are trained and combined with the frozen parameters. This includes approaches such as the use of Adapters (e.g., Pfeiffer Adapter) [51, 109], or optimising soft-prompt (e.g., prefix-tuning [71], p-tuning [79] or prompt-tuning [65]). Second, selective methods, do not introduce additional parameters, but instead choose a subset of the pre-trained model parameters to train. This includes approaches based on lottery ticket hypothesis [9] or other approaches for choosing the subsets [40, 156]. Third, reparametrised approaches, introduce additional low-rank trainable parameters during training that are then fused with the original model for inference. The most popular approach is the low-rank adaptation (LoRA) [52] and its alternatives. Finally, there are also hybrid approaches that combine multiple PEFT methods.

For a more comprehensive overview of different *parameter-efficient fine-tuning* approaches, please refer to [47].

C Implementation of the Survey Methodology

C.1 Search query definition

The search query is composed of two parts: 1) *terms* that represent all the versions of stability that we consider in this survey; and 2) *scope* describing our focus on approaches that can be used for learning with limited labelled data.

The set of terms that can be used to describe the effects of randomness on the performance was progressively expanded based on the terms in papers that in any way address this randomness. The identification of the papers was done randomly by first sampling papers that reference *stability*, assessing whether their understanding of stability adheres to our definition, and using the cited and citing papers to identify the next potential candidates. Using this methodology we identified the first set of key-words that included the following words: *stability*, *instability*, *sensitivity*, *variance*, *randomness*, *robustness* and *reproducibility*. However, we decided to remove the terms *robustness* and *reproducibility* from this list, as they have much wider meaning and would lead to many irrelevant papers.

The *scope* was defined based on the groups of approaches we want to focus on. Based on reading the papers, we have noticed only two terms being widely used to define all the approaches we were interested in: *meta-learning* and *few-shot learning*. We decided to use the term *few-shot learning* instead of *learning with limited labelled data*, because of its popularity and its misuse leading to both terms meaning the same thing in the literature. We noticed many different uses of these terms across papers, with many removing the hyphenation or foregoing the "learning" part of *few-shot learning*. To account for all the different uses of the terms, we finally decided to use the most simplified version of them, specifically *meta learning*, *metalearning* and *few shot*.

The final query we used for the keyword search was: ‘(“few shot” OR “meta learning” OR “metalearning”) AND (“stability” OR “instability” OR “sensitivity” OR “variance” OR “randomness”)’. We also searched only for papers that were published in the year 2017 or later, as the focus on the effects of randomness and stability before this year was non-existent.

C.2 Relevant Digital Libraries Used to Discover Papers

We used digital libraries that provide flexible search options, with a focus on the machine learning domain as one of its domains and allow for reproducible search results to discover the relevant papers. As we were not interested in other disciplines, we limited the search results to only those

relevant for us. This includes the following databases: ACM Digital Library, IEEE Xplore and Scopus database.

In addition, we also searched through the proceedings of the top-rated machine learning conferences, such as ICML, ICLR, or NeurIPS. We have noticed that many relevant papers are submitted to these conferences, while at the same time only a fraction of them are indexed by the previously stated databases. We specifically searched through ACL Anthology (containing the aggregated proceedings of the ACL conferences), the NeurIPS conference and the Proceedings of Machine Learning Research (PMLR), which includes proceedings from conferences such as ICML or ICLR. Only the ACL anthology allows for advanced search through papers, nevertheless, the search is limited to only the top 100 results. Therefore, to search through these libraries we opted to use Google Scholar. Each conference was searched separately using a modified query to limit results to only specific site (using the keyword in search, i.e. ‘site:aclanthology.org’ for ACL, ‘site:neurips.cc’ for NeurIPS and ‘site:mlr.press’ for PMLR). We take only the top 200 relevant papers from each conference. To determine the cut-off point, we manually determined the relevance of the papers at later pages of results and noticed the relevance drops of significantly after ~100 papers. Therefore we decided to search through 100 more papers just to be sure we were not missing any relevant papers.

Lastly, we also use included cited and citing papers from the most relevant ones as additional sources. We detected a significant clustering behaviour in papers that deal with the effects of randomness, with papers related to specific groups of approaches extensively citing all other papers in this group. This significantly improves the coverage of results as it includes papers not indexed in the databases, but also those that may use slightly different terminology or those that are relevant, but do not focus only on the limited data setting.

C.3 Identification of relevant papers using search and filtering

Since the research paper collection from digital libraries would naturally result in a large set of papers with various relevance, we selected the subset of the most relevant papers that we call *core* papers. To identify the set of these papers, we perform the following steps across multiple iterations:

- (1) **Create a set of potentially relevant papers.** The defined search query is used across all specified digital libraries to create the set of papers that are potentially relevant for the work. This generates a large set of papers with many false positives.
- (2) **Filter papers based on relevance.** The large set of potential papers is reduced to a set of papers dealing to some extent with the effects of randomness. To filter out the irrelevant papers the context in which the individual key-words from the search query are used in the paper. This is done by finding all the appearances of the key-words using a full-text search and manually checking them. After finding the first relevant context for both the *term* and *scope* key-words, the paper is considered to be relevant. If no such context exists for one of the groups, or the context indicates a different definition of the problem (such as *algorithmic stability*), the paper is considered to be irrelevant and removed from the set. For example, we do not consider the context to be relevant if it appears in related work. This step significantly reduces the set, leaving only 5% of papers.
- (3) **Filter papers based on merit.** Although the papers in this step are mostly relevant, the set still contains papers we do not want to consider based on their form. This includes papers that are only extended abstracts, proposals for talks or projects, or papers appearing on papers not relevant to stability. In addition, papers from conferences with a lower Core rank than B or from journals not belonging to Q2 or Q1 are also removed.

- (4) **Filter papers based on focus.** In this step we filter out the papers that only recognise that there is a problem with stability in regards to randomness in some approaches but provide no further investigation or focus on the randomness. We call these papers as *recognise* papers. This mostly includes papers that note that the training process with limited data is unstable so they perform multiple runs, or papers that note that it was found in previous papers that there is a problem. This leaves us with a set of papers we consider to be *core*.
- (5) **Extend the set of papers with citing and cited papers.** As there is a strong clustering present, papers dealing with similar approaches cite each other extensively. We make use of this clustering to discover other relevant papers that may not be indexed by the databases or that use different terminology so the search query does not catch them. However, we extract such papers only based on the context they are referenced in (when cited in the paper), or based on the relevant filtering step applied only to the title and abstract of the paper (which cites the current paper). We identify the citing papers through the Google Scholar *cited by* section present next to each paper.

We repeat the steps 2 – 5 until no new papers are identified in step 5.

C.4 Analysis and categorisation of the papers

All the papers identified this way were then manually categorised using our taxonomy as described in the main content of the survey (Section 3). The content of each paper was analysed and assigned value for each defined property. Afterwards, the main categorisation property (task performed to address the randomness) was selected based on how best it splits the papers into independent groups.

For each task for addressing the effects of randomness, additional properties and characteristics were identified and used to cluster the papers further. These additional properties were then used to derive the main findings from each task. Therefore, each finding can be mapped to one or multiple such properties at the same time. Besides the main properties and characteristics described in Section 3, we consider also the following additional properties and characteristics of the papers:

- (1) *Modality* that divides the papers into those focusing on text data, images, tabular data, or any combination of these modalities.
- (2) *Groups of approaches* further divide the primary property of the *machine learning approach* into a more detailed categorisation. For example, the *meta-learning* approaches are further divided into optimisation or metric-based meta-learning.
- (3) *Interactions*, a binary property that specifies whether the confounding effects between different *randomness factors* are addressed in any way in the paper.
- (4) *Standard deviation*, a binary property that specifies whether the performance is reported using standard deviation (or similar metric such as confidence intervals) or just as a simple aggregate metric (mean, median, etc.).
- (5) *Out-of-distribution*, a binary property that specifies whether the data used for testing purposes in the papers comes from the same distribution as the training data, such as addressing the effects of randomness when dealing with multilingual data or data drifts.
- (6) *Datasets* used to better explore what datasets are popular when addressing the effects of randomness.

The challenges and open problems were identified by observing common patterns in the findings, such as contradictions in results. In addition, the open problems were further extended by identifying parts of the research process that are currently missing even though their inclusion is a straightforward and simple modification of the current approach, such as using statistical evaluation of results instead of reporting a single performance metric.