

# Invariant Language Modeling

Anonymous ACL submission

## Abstract

Modern pretrained language models are critical components of NLP pipelines. Yet, they suffer from spurious correlations, poor out-of-domain generalization, and biases. Inspired by recent progress in causal machine learning, in particular the invariant risk minimization (IRM) paradigm, we propose *invariant language modeling*, a framework for learning invariant representations that generalize better across multiple environments. In particular, we adapt a game-theoretic implementation of IRM (*IRM-games*) to language models, where the invariance emerges from a specific training schedule in which all the environments compete to optimize their own environment-specific loss by updating subsets of the model in a round-robin fashion. In a series of controlled experiments, we demonstrate the ability of our method to (i) remove structured noise, (ii) ignore specific spurious correlations without affecting global performance, and (iii) achieve better out-of-domain generalization. These benefits come with a negligible computational overhead compared to standard training, do not require changing the local loss, and can be applied to any language model architecture. We believe this framework is promising to help mitigate spurious correlations and biases in language models.

## 1 Introduction

Despite dramatic progress in NLP tasks obtained by modern pretrained transformer models, important limitations remain. In particular, pretrained language models suffer from poor generalization, even under small perturbations of the input distribution (Moradi and Samwald, 2021). Indeed, these models encode (Moradi and Samwald, 2021) and exploit (Tu et al., 2020; Niven and Kao, 2019) spurious correlations, i.e., correlations that do not generalize across data distributions. Since language models are trained on large unverified corpora, they also suffer from biases (Nadeem et al., 2021; Bordia and Bowman, 2019). Biases are correlations

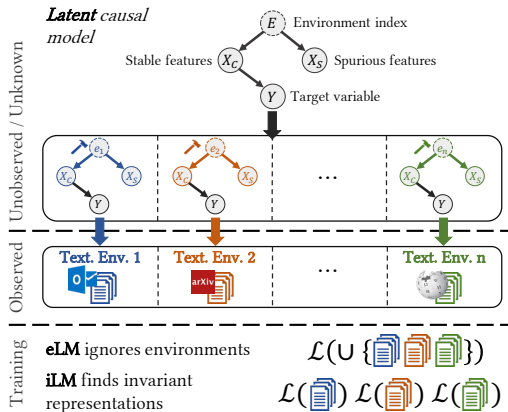


Figure 1: **High-level overview using a simplified causal structure.** The distinction between environments makes it possible to separate spurious from stable features. Indeed, the relationship between the target variable  $Y$  and the stable features  $X_C$  is invariant across environments:  $\mathbb{E}[Y|X_C, E] = \mathbb{E}[Y|X_C]$ . However, the correlation between  $Y$  and  $X_S$  is spurious and does not generalize across environments:  $\mathbb{E}[Y|X_S, E = e] \neq \mathbb{E}[Y|X_S, E = e'], e \neq e'$ . Language models trained with the standard ERM, denoted as eLM in this work, exploit all correlations available during training and aim to learn  $\mathbb{E}[Y|X_C, X_S]$ . Our proposed invariant language models, denoted as iLM, focus on invariant features and aim to learn  $\mathbb{E}[Y|X_C]$ . In language modeling,  $Y$  could represent the missing-word prediction task.

that may or may not be spurious according to the available textual data distributions but are nevertheless undesired. Existing techniques aiming to remove spuriousness or biases involve computationally expensive domain alignment (Akuzawa et al., 2019; Liu et al., 2020; Zhao et al., 2020), domain transfer (Balaji et al., 2018) or adding penalty terms in the loss targeted at specific undesired correlations (Qian et al., 2019; Zhao et al., 2018). Alternatively, data preprocessing (Zhao et al., 2017; Zhou et al., 2021) or manipulation such as counterfactual data-augmentation (Lu et al., 2018) can yield datasets where the undesired correlations are

056 less present. Pretraining with larger and more di- 107  
057 verse datasets can also help (Tu et al., 2020; Brown 108  
058 et al., 2020).

059 However, recent works on the theory of causality 109  
060 (Pearl, 2018; Schölkopf, 2019) argue that removal 110  
061 of spurious correlations requires altogether differ- 111  
062 ent learning and training paradigms going beyond 112  
063 purely statistical learning. Indeed, generalization, 113  
064 spuriousness, and biases are all better understood in 114  
065 the language of causality (Pearl, 2018). Intuitively, 115  
066 causal relationships are the ones expected to be sta- 116  
067 ble (Schölkopf et al., 2021; Peters et al., 2017) and 117  
068 generalizable (Peters et al., 2016). When the causal 118  
069 graph underlying the data generation mechanism is 119  
070 known, there exist causal identification algorithms 120  
071 to distinguish *desired* from *undesired* correlations 121  
072 (Shpitser and Pearl, 2008). However, for complex 122  
073 tasks of interest, the underlying causal model is not 123  
074 known. Language modeling is one of these tasks, 124  
075 where it is unclear what would even be the relevant 125  
076 random variables constituting the causal model. 126

077 Therefore, causal identification from the causal 127  
078 graph seems out-of-reach for language modeling. 128  
079 Similarly, removing undesired correlations one by 129  
080 one is impractical due to the sheer amount of possi- 130  
081 ble correlations to consider. In this work, we 131  
082 propose to leverage recent progress in causal ma- 132  
083 chine learning to offer a new and more flexible 133  
084 lever for dealing with spuriousness and biases. 134  
085 We take inspiration from the *invariance princi- 135  
086 ple*, which states that only relationships invariant 136  
087 across training *environments* should be learned (Pe- 137  
088 ters et al., 2016). Under specific assumptions, the 138  
089 invariant representation would then only encode 139  
090 the causal relationships relevant to the task and 140  
091 should thus generalize. Environments correspond 141  
092 to different views of the learning task, i.e., dif- 142  
093 ferent data distributions. The invariance princi- 143  
094 ple is illustrated by Fig. 1 with a simplified causal 144  
095 model as an example.  $E$  represents environment 145  
096 indices,  $Y$  is the target variable,  $X_C$  are the *causal 146  
097 features*, such that  $\mathbb{E}[Y|X_C]$  is stable across envi- 147  
098 ronments ( $\mathbb{E}[Y|X_C, E] = \mathbb{E}[Y|X_C]$ ), and  $X_S$  are the 148  
099 spurious features, not generalizing across environ- 149  
100 nments ( $\mathbb{E}[Y|X_S, E = e] \neq \mathbb{E}[Y|X_S, E = e'], e \neq e'$ ). 150  
101 Language models trained with standard empirical 151  
102 risk minimization (ERM), denoted as eLM in this 152  
103 work, exploit all correlations available during train- 153  
104 ing and aim to learn  $\mathbb{E}[Y|X_C, X_S]$ . Our proposed 154  
105 invariant language models, denoted as iLM, focus 155  
106 on invariant features and aim to learn  $\mathbb{E}[Y|X_C]$ . In 156  
157

107 practice, since the causal model is unknown, it 108  
109 is the choice of environments that defines what 110  
111 correlations are spurious. Invariant learning with 112  
113 appropriate choices of environments is the lever 114  
115 we propose to employ to more flexibly deal with 116  
117 spuriousness and biases. 118

119 A practical implementation of the invariance 120  
121 principle was proposed by Arjovsky et al. (2019). 122  
123 They introduced *invariant risk minimization* (IRM), 124  
125 an alternative to ERM as a training objective enforc- 126  
127 ing the learning of invariant representations. Ahuja 128  
129 et al. (2020) later improved the training procedure 130  
131 to solve the IRM objective with a method called 131  
132 IRM-games. Unlike previous methods for remov- 133  
134 ing biases and spurious correlations, IRM-games 134  
135 does not modify the loss with a regularization 135  
136 term and does not compute domain alignment (or 136  
137 matching) statistics. The invariance benefits come 137  
138 from the specific training schedule where environ- 138  
139 ments compete to optimize their own environ- 139  
140 ment-specific loss by updating subsets of the model in 140  
141 a round-robin fashion. The Nash equilibrium of 141  
142 this game between environments is a solution to 142  
143 the IRM objective (Ahuja et al., 2020). 143

144 We argue that the IRM paradigm, and IRM- 144  
145 games specifically, is well-suited to improve mod- 145  
146 ern NLP systems. Textual data naturally comes 146  
147 from different environments, e.g., encyclopedic 147  
148 texts, Twitter, news articles, etc. Moreover, not 148  
149 knowing the causal mechanisms behind language 149  
150 generation within these environments is not a 150  
151 blocker, as the relevant variables can now remain 151  
152 latent. In this work, we adapt IRM-games to lan- 152  
153 guage modeling. This involves continuing the train- 153  
154 ing of existing pretrained models to enforce invari- 154  
155 ant representations. We then investigate the ability 155  
156 of iLM to remove undesired correlations in a series 156  
157 of controlled experiments, effectively answering 157  
158 our core **research question**: Does the invariance 158  
159 principle give rise to a practical strategy to remove 159  
160 undesired correlations from language models? 160

161 **Contributions.** (i) We introduce a new training 161  
162 paradigm (iLM) for language models based on the 162  
163 invariance principle (Sec. 3). Thanks to the use of 163  
164 the IRM-games training schedule (see Sec. 2), our 164  
165 iLM framework results in negligible computational 165  
166 overhead compared to standard ERM training, does 166  
167 not require changing the local loss, and is agnostic 167  
168 to the language model architecture. (ii) In a se- 168  
169 ries of controlled experiments (Sec. 4), we demon- 169  
170 strate the ability of iLM to remove structured noise 170  
171

(Sec. 4.1), ignore specific spurious correlations without affecting global performance (Sec. 4.2), and achieve better out-of-domain generalization (Sec. 4.3). (iii) We discuss our contributions in relation to previous work (Sec. 5). (iv) Finally, we release Huggingface-compatible code for training iLM using existing language model checkpoints (Wolf et al., 2020): [anonymized](#)

## 2 Background

In this section, we present the ideas and previous work necessary to understand our proposed models.

### 2.1 Invariance across Environments (IaE)

Recent works on the theory of causality (Pearl, 2018; Schölkopf, 2019) have argued that out-of-distribution generalization and removal of spurious correlations require going beyond purely statistical learning. This is motivated by the intuition that causal relationships are the ones that are expected to be robust and generalizable (Peters et al., 2016). Unfortunately, for problems of interest, the causal model is usually unknown. Then, different methods based on different assumptions can still hope to capture some causal properties important for generalization, e.g., ensuring that only causal parents of the target variable are used for prediction. In causal machine learning, these ideas crystallized in the *invariance principle* which states that only relationships invariant across training environments should be learned (Peters et al., 2016; Muandet et al., 2013). In this paradigm, different environments correspond to data collected in different setups, i.e., different data distributions (Pearl, 2018). Interestingly, learning according to the invariance principle does not require knowing what modifications of the data generation mechanism happened in which environment, it only requires that  $\mathbb{E}[Y|X_C]$  remains unchanged, where  $X_C$  are the causal parents of the target variable  $Y$  (Arjovsky et al., 2019; Rosenfeld et al., 2021).

### 2.2 Invariant Risk Minimization (IRM)

While the invariance principle is a general and powerful idea, works based on this principle often require knowing which random variables are part of the causal model (Akuzawa et al., 2019; Peters et al., 2016). Arjovsky et al. introduced *invariant risk minimization* (IRM), an alternative to empirical risk minimization, and a practical training objective compliant with the invariance principle. IRM

allows for relevant variables to remain latent. Under specific assumptions, it will ignore correlations not due to the causal parents of the target variables.

IRM builds on the idea that the training data comes from different environments  $e \in \mathcal{E}$ . Each environment  $e \in \mathcal{E}$  induces i.i.d. samples  $D^e$  from a distribution  $P(X^e, Y^e)$ . Then, the goal is to use these multiple datasets to learn a predictor  $Y \approx f(X)$ , which performs well across the set of all environments  $\mathcal{E}^*$ , only part of which were seen during training:  $\mathcal{E} \subset \mathcal{E}^*$ . This is accomplished by decomposing  $f$  into a feature representation  $\phi$  and a classifier  $w$  as  $f = w \circ \phi$ , where  $\circ$  denotes function composition, i.e.,  $(w \circ \phi)(X) = w(\phi(X))$ . The feature representation  $\phi$  elicits invariant representation of the data if the same classifier  $w$  is simultaneously optimal for all environments  $e \in \mathcal{E}$ . Thus, IRM solves the following optimization problem:

$$\min_{\phi, w} \sum_{e \in \mathcal{E}} R^e(w \circ \phi), \quad (1)$$

$$\text{subject to } w \in \arg \min_{w'} R^e(w' \circ \phi), \forall e \in \mathcal{E}, \quad (2)$$

where  $R^e$  is the empirical risk computed within an environment  $e$ ; i.e., if  $\mathcal{L}$  is a loss function,  $R^e = \mathbb{E}[\mathcal{L}((w \circ \phi)(X^e), Y^e)]$ .

### 2.3 IRM-games

IRM is a challenging bi-level optimization originally solved (Arjovsky et al., 2019) by relaxing the objective function, setting the invariance criteria as a regularizer. Later, Ahuja et al. improved the training procedure by using a game-theoretic perspective in which each environment  $e$  is tied to its own classifier  $w^e$ , and the feature representation  $\phi$  is shared. The global classifier  $w$  is then defined as the ensemble of all environment-specific classifiers. Environments take turns to make a stochastic gradient update to minimize their own empirical risk  $R^e(w \circ \phi)$  but the update concerns only their own classifier  $w^e$ , while the shared  $\phi$  is updated periodically. For more details see the algorithm called V-IRM in the original paper. Ahuja et al. showed that the equilibrium of this game is a solution to the IRM objective.

## 3 Model

We introduce a way to train language models inspired from the IRM-games setup. This involves distinguishing the shared invariant feature learner

$\phi$  from the environment specific  $w_e$ 's. With modern language models architectures, a natural choice emerges:  $\phi$  as the main body of the encoder, and  $w_e$  as the language modeling head that outputs the logits after the last layer.

Formally, suppose we have  $n$  environments consisting of data  $\{(X^e, Y^e)\}_{e=1 \dots n}$ . For a batch  $(x_i, y_i) \sim P(X^i, Y^i)$  from environment  $i$ , the model output is formed using an ensemble of  $n$  language modeling heads  $\{w_e\}_{e=1 \dots n}$  on top of the transformer encoder:  $\hat{y} = \text{softmax}\left(\frac{1}{n} \sum_{e=1}^n w_e \circ \phi(x_i)\right)$ . Then, a (masked) language modeling loss  $\mathcal{L}$  is computed on the model output  $\hat{y}$ . Note that it is the predictions of the  $n$  heads that are averaged (compared to the weights or gradients as in a multi-task setup). No head gets to predict alone; the  $n$  heads always predict together as an ensemble but performing competitive gradient updates in a round-robin fashion, which in turn creates the game-theoretic conditions that enforces the invariance of Eq. 1.

**Training** The training of iLM follows the pseudocode described in Alg. 1, where environments take turn to send a batch of data and update  $\phi$  and their associated head. An illustration is provided in Appendix A. Each head periodically gets an opportunity to pull the global ensemble classifier  $\mathbf{w}$  and the feature learner  $\phi$  towards fitting the distribution of its associated environment. Intuitively, since each head gets the same amount of updates, the game converges to a global classifier that is simultaneously optimal for each environment, as demonstrated by (Ahuja et al., 2020). If the model one head per environment trained in round-robin fashion but without the ensemble prediction and competitive gradient update (similar to multi-task learning), it would not enforce invariance across environments.

While the V-IRM algorithm of Ahuja et al. (2020) only updates  $\phi$  periodically, we found it more stable to update it together with every head update.

**Advantages of design choices** Choosing the heads as environment-specific  $w_e$  is agnostic to the model architecture because the whole body of the model is included in  $\phi$ . Only the components specific to language modeling – the heads – have a different structure compared to the standard ERM setup. This makes the iLM framework compatible with any kind of pretrained language model. Moreover, the whole body of the model is the in-

---

### Algorithm 1 iLM training

---

```

1: Initialize( $\phi$ )
2: Initialize( $\{w_e\}_{e \in \mathcal{E}}$ )
3: for iteration  $\in \{1, 2, \dots, \frac{N_{steps}}{|\mathcal{E}|}\}$  do
4:   for environment  $i \in \mathcal{E}$  do
5:      $(x_i, y_i) \leftarrow \text{GetBatchFromEnv}(e)$ 
6:     CompetitiveUpdate( $x_i, y_i, \phi, \{w_e\}_{e \in \mathcal{E}}$ )
7:   end for
8: end for
9: function COMPETITIVEUPDATE( $x_i, y_i, \phi, \{w_e\}_{e \in \mathcal{E}}$ )
10:   $L = \mathcal{L}\left(\text{softmax}\left(\frac{1}{n} \sum_{e=1}^n w_e \circ \phi(x_i)\right), y_i\right)$ 
11:  GradientUpdate( $L, \phi, w_i$ )
12: end function

```

---

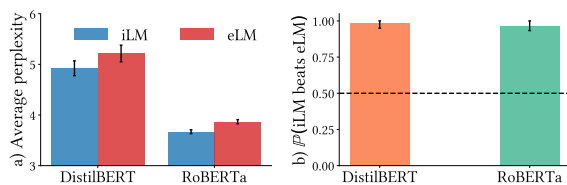


Figure 2: **Structured noise removal experiment:** a) average perplexity over all hyper-parameters, b) Probability that iLM has a lower perplexity than eLM when compared on the same hyper-parameters.

variant feature learner  $\phi$ . Finally, since only the heads and their training dynamic differ from standard eLM, the usage of iLM models does not differ in downstream tasks.

## 4 Experiments

Invariance training comes with the promise of robustness and generalization (Peters et al., 2016; Muandet et al., 2013; Ahuja et al., 2020). In the following series of experiments, we test whether our proposed architecture for language modeling can provide such benefits. Since the causal model governing language production is unknown, we do not have access to the gold standard answer about which correlation is spurious. Thus, we focus on controlled setups: crafting environments whose difference is known, from which we know the expected behavior. We describe three main experiments: structured noise removal, controlled correlation removal, and out-of-domain generalization. We emphasize that we use perplexity evaluation in two out of three experiments, not because we view low perplexities as desirable for language models, but because perplexity is an objective measure of the ability of a language model to fit data



that matches its training goal. Perplexity evaluation is part of the simplified and controlled setup used to test the new core benefits of iLM. The results presented here open the way for more practical future works based on what we call **environment design**: how to choose environment splits to be useful in downstream tasks (see Sec. 5 for an extended discussion).

For all the experiments, each plot reports 95% confidence intervals from bootstrap resampling of the data. We repeat each experiment for two base pretrained transformer models with different properties (size, tokenization method): distilBERT (Sanh et al., 2019) and ROBERTa (Liu et al., 2019). We also repeat each experiment with different learning rates, number of training steps and random restarts with different random seeds. Appendix B provides additional details regarding each experiment and further results about the importance of hyper-parameters.

### 4.1 Structured Noise Removal

**Description.** In this experiment, we test robustness in a controlled setup. We craft two environments: Env-A made of clean Wikipedia articles and Env-B made of full HTML pages of Wikipedia articles. Then, we continue the training with the masked language modeling (MLM) loss from existing checkpoints for both iLM and eLM with these two environments and evaluate the MLM perplexity on a held-out dataset of clean Wikipedia articles. Intuitively, eLM should try to fit the HTML part of the training data and thus be more surprised by the clean Wikipedia articles during the test set. However, iLM should learn to ignore the HTML because it does not generalize from Env-B to Env-A.

The results are visualized in Fig. 2. See Appendix B.1 for hyper-parameters considered. On the left plot, we report the average perplexity on the test set averaged over all experiments. On the right plot, we report the probability that for any given hyper-parameter configuration, iLM has a lower perplexity than eLM. In these experiments, paired comparison is particularly important because varying hyper-parameters results in large variations of perplexity. Blindly averaging amplifies the variance and hides the structure of model performance (Peyrard et al., 2021). For reference, the perplexities on the same test set of pretrained distilBERT and ROBERTa are, respectively, 14.43 and 6.71.

**Analysis.** We observe that iLM has an overall lower test perplexity when averaged over all experiments (Fig. 2 a). Furthermore, for any given hyper-parameter choice, iLM is better than eLM (Fig. 2 c) with a probability  $> .95$  for both distilBERT and ROBERTa. Note that the few cases where eLM matches or beats iLM happen when few training steps have been taken ( $< 50$ ). The trends are the same for both distilBERT and ROBERTa despite large perplexity differences between them.

### 4.2 Controlled Correlation Removal

**Description.** In this experiment, we test the capacity to remove one precise and known correlation by crafting two environments differing only in this specific correlation. We use binarized gendered terms and create two environments where the gendered terms are used differently.<sup>1</sup> More precisely, we take a textual data source with known gender bias, in this case, Wikitext-2 (Merity et al., 2016). A fraction  $p$  of the data goes into Env-A, the rest  $(1 - p)$  goes into Env-B. Env-A remains untouched and preserves all the properties of the original data source. Whereas Env-B is intervened upon by inverting all gendered terms based on a dictionary provided by previous work (Bordia and Bowman, 2019). When  $p = 1 - p = 0.5$ , this setup matches the counterfactual data-augmentation methods (Lu et al., 2018) already used to mitigate gender-bias in language models. Intuitively, iLM should learn to ignore gender-based correlations no matter what is the fraction  $p$ . However, eLM is only expected to ignore them when  $p = 1 - p = 0.5$ , i.e., the two environments have the same number of samples (Lu et al., 2018).

We craft this experiment as an example of controlled correlation removal, but it shows promise for practical bias removal because selecting or crafting environments where biases do not hold is arguably simpler than precisely counter-balance the bias by data processing/augmentation or regularization. iLM can directly improve current bias-removal strategies based on counterfactual data augmentation. We come back to this in Sec. 5.

**Experimental setup.** To measure whether the correlation has been successfully removed: (i) we take

<sup>1</sup>We recognize the non-binary nature of gender as well as the many ethical principles in the design, evaluation, and reporting of results in studying gender as a variable in NLP (Larson, 2017). Because iLM is not limited to training only with two environments, this architecture can also support more general bias removal goals.

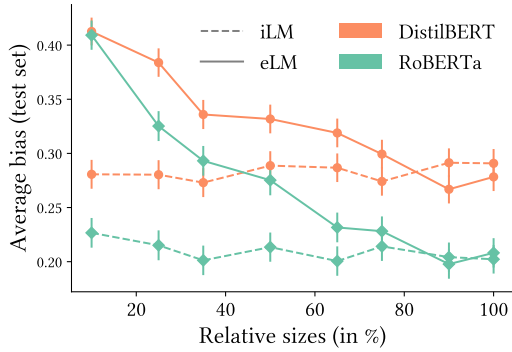


Figure 3: **Controlled correlation removal experiment:** On the x-axis, we report the relative size between the modified environment and the unmodified one, and report on the y-axis the average bias for both iLM and eLM. Note that  $\mathbb{P}(iLM \text{ beats } eLM) > 0.95$  when the relative size is  $< 80\%$ , eLM and iLM become indistinguishable for relative sizes  $> 80\%$ .

all gendered terms in the test set, (ii) replace them by the MASK token, (iii) use trained eLM and iLM models to predict the missing term, (iv) look in the softmaxes the scores received by the terms of the target gendered-pair. We note  $s_f$  and  $s_m$  the score assigned to the female and male terms in the softmax. (v) Finally, we compute an entropy-bias measure:  $B_H = H_2\left(\frac{1}{2}\right) - H_2\left(\frac{s_f}{s_f + s_m}\right)$ , where  $H_2$  is the binary entropy (note that  $H_2\left(\frac{1}{2}\right) = 1$ ).  $B_H$  measures the extent to which a softmax has a preference for the male or female term in a gendered pair of terms. For example, in the sentence "MASK is the best doctor", we look at the softmax score of the gendered-pair [he, she]. If a model has learned to ignore gender-based correlation, the entropy should be high, i.e., which gender to be used is uncertain and the entropy bias  $B_H$  should be low.

We ran the experiments for varying values of  $p$  and report the results in Fig. 3. See Appendix B.2 for hyper-parameters considered. For reference, the entropy bias of distilBERT and RoBERTa before training are, respectively, 0.39 and 0.46.

**Analysis.** Both eLM and iLM decrease the average entropy bias in the balanced setup but only iLM succeeds in the unbalanced setup. In the balanced setup (relative sizes close to 100%), eLM and iLM perform within each other’s confidence intervals. However, in the unbalanced setup, iLM largely outperforms eLM. We note that the probability that iLM beats eLM for any given hyper-parameter configuration is  $> 0.9$  for both distilBERT and RoBERTa when the relative sizes is below 80%. As

desired iLM is not affected by the relative size of the environments. These results confirm the hypothesis, that correlation reduction needs a precisely balanced dataset for eLM (Lu et al., 2018), while it matter much less for iLM. Furthermore, this entropy bias reduction does not happen at the cost of worst general perplexities (See Appendix B.2).

### 4.3 Out-of-domain Generalization

In this experiment, we venture beyond carefully controlled setups and test out-of-domain generalization with naturally occurring domains. We use subsamples from *thePile* dataset (Gao et al., 2020) which contains 20 very diverse textual domains: OpenSubtitles, ArXiv papers, News, GitHub comments, etc.

**Experimental setup.** We randomly sample  $n$  domains from thePile, use  $n - 1$  of these domains as training and the remaining unseen one for testing. We compare iLM and eLM on their ability to generalize on the unseen domains by measuring the perplexity on the test domain.

The disparity of domains in thePile results in vast differences in perplexities between different domains, making the perplexities not comparable from one testing domains to the next. Instead of reporting averages of different domains, we report the better suited paired evaluation: comparing iLM and eLM in the same experimental setup (same hyper-parameters and same training/testing domains). The probability that iLM is better than eLM after 5000 training step is 0.9 with the 95% confidence interval of (0.79, 1). In Appendix B.3, we provide details about the impact of hyper-parameters.

However, the advantage of iLM over eLM is less striking in this experiment than in the two previous ones. The average perplexities of iLM is not always significantly lower than that of eLM (see Appendix B.3 for details). We come back to potential reasons for this behavior in Sec. 5.

## 5 Discussion

In this section, we discuss our contributions in the context of previous work.

### 5.1 Related Work

**Domain generalization.** The performance of deep learning models substantially degrades on Out-of-Domain (OoD) datasets, even in the face of small variations of the data generating process (Hendrycks and Dietterich, 2019). Blanchard et al.

(2011) have proposed domain generalization (DG) as a formalism for studying this problem. In DG, the goal is to learn a model using data from a single or multiple related but distinct training domains, in such a way that the model generalizes well to any OoD testing domain, unknown during training. Recently, the problem of DG has attracted a lot of attention, and has been approached from different facets. Most of the existing methods fall under the paradigm of domain alignment (Muandet et al., 2013; Li et al., 2018b; Akuzawa et al., 2019; Liu et al., 2020; Zhao et al., 2020). Motivated by the idea that features that are stable across the training domains should also be robust to the unseen testing domains, these methods try to learn domain-invariant representations. A group of other methods is based on meta-learning (Dou et al., 2019; Balaji et al., 2018; Li et al., 2018a). The motivation behind this approach is that it exposes the model to domain shifts during training, which will allow it to generalize better during testing. Regularization through data augmentation is commonly used in the training of machine learning models to alleviate overfitting and thereby improve generalization (Zhou et al., 2021, 2020).

**Domain generalization applied to language models.** In NLP, the default pipeline involves pre-training a task-agnostic language model, which is then finetuned on downstream tasks. This pre-training/finetuning division of learning is already known to improve robustness on downstream tasks (Hendrycks and Dietterich, 2019). However, the language models themselves suffer from spurious correlations and poor generalization even with small perturbations of the inputs (Moradi and Samwald, 2021). To alleviate such problems, Oren et al. (2019) adapt Distribution Robust Optimization (Ben-Tal et al., 2013) to language models. This results in a new loss minimizing the worst-case performance over subsamples of the training set. They focus on domains with topic shifts. Then, Vernikos et al. (2020) use domain adversarial regularization to improve testing performance on unseen domains.

Also related to our framework are techniques aiming at de-biasing language models. Biases are correlations that may or may not be spurious but are nevertheless undesired. Removing such biases is typically done by (i) adding a bias-specific penalty term (Qian et al., 2019; Bordia and Bowman, 2019; Zhao et al., 2018) to the loss, and/or (ii) augmenting the data to counterbalance the undesired correlation

(Lu et al., 2018; Zhao et al., 2017). For example, counterfactual data-augmentation used to reduce gender-bias (Lu et al., 2018) flips half of the gendered terms to destroy existing correlations in the original inputs.

**Justification of IRM-games.** The rich literature in domain generalization begets the question why we should focus specifically on IRM-games to adapt to language models. Counterfactual data augmentation techniques require some knowledge of and some ability to manipulate the possible mechanisms generating the data. Meta-learning techniques come with a large extra-computation cost as they are based on multiple rounds of training. This is not practical for modern language models. IRM-games lends itself particularly well to modern implementations of language models with the natural distinction between the transformer body as  $\phi$  and the language modeling heads as  $w$ . Importantly, as opposed to most other methods, it does not require extra computation about the environments (like matching, variance, drift, etc.). It is sufficient to keep track of environment indices during training and the invariance comes from the particular game-theoretic dynamics of the training schedule. Thus, the local language modeling loss can remain unchanged, there is no need for a regularization term for which the strength needs to be tuned. Finally, iLM has a minimal computational overhead compared to eLM because only the heads are multiplied (one per environment) but the number of parameters in these heads is small in comparison to the number of parameters in the main body a modern language model.

## 5.2 Potential Limitations of Domain Generalization Methods

**Discussion of potential limitations.** With the recent interest in invariance-based methods came a other works questioning the real generalization ability of these methods. For example, Gulrajani and Lopez-Paz (2021) finds that finetuning ERM can be as good as vanilla IRM (Arjovsky et al., 2019). Similarly, Rosenfeld et al. (2021) find that the number of environments needed for *full* generalization can be *large*. To organize the discussion around the benefits of OoD generalization methods, Ye et al. (2021) argue about the importance of distinguishing different types of distribution shifts according to the underlying data generation mechanism. In particular, they distinguish *diversity shifts*



and *correlation shifts*, and claim that invariance-based methods perform well for correlation shifts but not for diversity shifts.

**In the language context.** These limitations did not include IRM-games as part of their analysis. In language, since the latent causal model is unknown, it is difficult to anticipate which kind of distribution shifts our models might face. Nevertheless, the experiments of structured noise removal (Sec. 4.1) and controlled correlation removal (Sec. 4.2) are instances of correlation shifts as defined by Ye et al. (2021). On these experiments, we observe striking improvements when compared to eLM. The OoD experiment (Sec. 4.3) involves more latent variables in the shifts from one domain to another and *possibly* exhibits both correlation and distribution shifts. This can explain the smaller performance gains observed in this experiment.

**Possible problems with environment choices.** One question that might arise from the iLM training schedule is what happens when environments have no lexical overlap? Maybe no correlation remains in iLM? To demonstrate that iLM operates on latent variables and not just on surface-level correlations, we perform a simple experiment with languages as environments. We train iLM with a pretrained multilingual model (XLM-ROBERTa) using English Wikipedia articles and Farsi Wikipedia articles as two environments. Despite absolutely no surface-level overlap, iLM is still able to improve perplexity in each language individually and does not destroy previously learned correlations. This experiment is detailed in Appendix B.4.

Also, if the number of environments grows arbitrarily large, certainly iLM would not find any stable correlations in the data. However, the choice of environments is not intended to be arbitrary; throwing as many environments as possible could not be expected to be useful. The choice of environments has to reflect assumptions about the underlying data generation mechanism. iLM then leverages the assumptions encoded in the choice of environments.

### 5.3 Environment Design

**Causal perspective.** Pearl organized causal problems in a three-level hierarchy termed the “ladder of causation”: *observational* queries correspond to seeing and observing; *interventional* queries correspond to acting and intervening; and *counterfactual*

queries correspond to imagining, reasoning, and understanding. In this ladder, it is in general impossible to solve problems at the higher-levels with only data and assumptions from lower levels. When performing invariant feature learning, we hope for generalization benefits from the interventional level (Peters et al., 2016; Arjovsky et al., 2019; Ahuja et al., 2020). However, ERM training and eLM operate at the observational level. The iLM setup also uses only observational data because the model is not performing experiments. Therefore, we need to inject causal assumptions (interventional level) to hope to get generalization benefits. These assumptions are encoded by the choice of environments (Peters et al., 2016; Arjovsky et al., 2019), which dictates where the interventions have happened in the unobserved data-generating process.

**Environment design.** This work has shown that iLM can effectively remove unstable correlations, the next question becomes that of **environment design**: *how to choose environment splits to be useful in practice?* or equivalently, *what assumptions are useful for tasks of interest?* Useful environment splits will likely be different for different tasks and different purposes. This work already demonstrated that the new paradigm of (i) environment design then (ii) iLM is practical for language-related problems. Simple environment choices already improve robustness, subsumes existing bias removal strategies, and are useful for OoD generalization. Choosing environment splits is a flexible way to inject priors and assumptions compared to manually deciding which correlation are desired (as in bias removal) or fully learning the causal graph (as in causal reasoning).

## 6 Conclusion

We introduce invariant language models trained to learn invariant feature representations that generalize across different training environments. In a series of controlled experiments, we demonstrate the ability of our method to remove structured noise, ignore specific spurious correlations without affecting global performance, and perform better out-of-domain generalization. These benefits come with a negligible computational overhead compared to standard training, do not require changing the loss, and apply to any language model architecture. We believe this framework is promising to help alleviate the reliance on spurious correlations and the presence of biases in language models.



701  
702  
703  
704  
705  
706  
707  
708  
  
709  
710  
711  
712  
713  
714  
  
715  
716  
717  
  
718  
719  
720  
721  
722  
723  
724  
  
725  
726  
727  
728  
729  
  
730  
731  
732  
733  
734  
735  
736  
  
737  
738  
739  
740  
741  
742  
743  
  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
  
755  
756

## References

Kartik Ahuja, Karthikeyan Shanmugam, Kush R. Varshney, and Amit Dhurandhar. 2020. [Invariant risk minimization games](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 145–155. PMLR.

Kei Akuzawa, Yusuke Iwasawa, and Yutaka Matsuo. 2019. [Adversarial Invariant Feature Learning with Accuracy Constraint for Domain Generalization](#). In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 315–331. Springer International Publishing.

Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. 2019. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*.

Yogesh Balaji, Swami Sankaranarayanan, and Rama Chellappa. 2018. [Metareg: Towards domain generalization using meta-regularization](#). In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 1006–1016.

Aharon Ben-Tal, Dick den Hertog, Anja De Waegenare, Bertrand Melenberg, and Gijs Rennen. 2013. [Robust solutions of optimization problems affected by uncertain probabilities](#). *Management Science*, 59(2):341–357.

Gilles Blanchard, Gyemin Lee, and Clayton Scott. 2011. [Generalizing from several related classification tasks to a new unlabeled sample](#). In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain*, pages 2178–2186.

Shikha Bordia and Samuel R. Bowman. 2019. [Identifying and reducing gender bias in word-level language models](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 7–15, Minneapolis, Minnesota. Association for Computational Linguistics.

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

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of](#)

[deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. 2019. [Domain generalization via model-agnostic learning of semantic features](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 6447–6458.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. [The pile: An 800gb dataset of diverse text for language modeling](#). *arXiv preprint arXiv:2101.00027*.

Ishaan Gulrajani and David Lopez-Paz. 2021. [In search of lost domain generalization](#). In *International Conference on Learning Representations*.

Ruocheng Guo, Pengchuan Zhang, Hao Liu, and Emre Kiciman. 2021. [Out-of-distribution prediction with invariant risk minimization: The limitation and an effective fix](#). *CoRR*, abs/2101.07732.

Dan Hendrycks and Thomas G. Dietterich. 2019. [Benchmarking neural network robustness to common corruptions and perturbations](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Brian Larson. 2017. [Gender as a variable in natural-language processing: Ethical considerations](#). In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 1–11, Valencia, Spain. Association for Computational Linguistics.

Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M. Hospedales. 2018a. [Learning to generalize: Meta-learning for domain generalization](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 3490–3497. AAAI Press.

Ya Li, Mingming Gong, Xinmei Tian, Tongliang Liu, and Dacheng Tao. 2018b. [Domain generalization via conditional invariant representations](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, pages 3579–3587. AAAI Press.

Chang Liu, Xinwei Sun, Jindong Wang, Tao Li, Tao Qin, Wei Chen, and Tie-Yan Liu. 2020. [Learning causal semantic representation for out-of-distribution prediction](#). *CoRR*, abs/2011.01681.

813	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. <a href="#">Roberta: A robustly optimized BERT pretraining approach</a> . <i>CoRR</i> , abs/1907.11692.	868
814		869
815		870
816		871
817		872
818	Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2018. <a href="#">Gender bias in neural natural language processing</a> . <i>CoRR</i> , abs/1807.11714.	873
819		874
820		875
821		
822	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. <a href="#">Pointer sentinel mixture models</a> . <i>CoRR</i> , abs/1609.07843.	
823		
824		
825	Milad Moradi and Matthias Samwald. 2021. <a href="#">Evaluating the robustness of neural language models to input perturbations</a> . <i>CoRR</i> , abs/2108.12237.	
826		
827		
828	Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. 2013. <a href="#">Domain generalization via invariant feature representation</a> . In <i>Proceedings of the 30th International Conference on Machine Learning</i> , volume 28 of <i>Proceedings of Machine Learning Research</i> , pages 10–18, Atlanta, Georgia, USA. PMLR.	
829		
830		
831		
832		
833		
834	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. <a href="#">StereoSet: Measuring stereotypical bias in pretrained language models</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 5356–5371, Online. Association for Computational Linguistics.	
835		
836		
837		
838		
839		
840		
841		
842	Timothy Niven and Hung-Yu Kao. 2019. <a href="#">Probing neural network comprehension of natural language arguments</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4658–4664, Florence, Italy. Association for Computational Linguistics.	
843		
844		
845		
846		
847		
848	Yonatan Oren, Shiori Sagawa, Tatsunori B. Hashimoto, and Percy Liang. 2019. <a href="#">Distributionally robust language modeling</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4227–4237, Hong Kong, China. Association for Computational Linguistics.	
849		
850		
851		
852		
853		
854		
855		
856	Judea Pearl. 2018. <a href="#">Theoretical impediments to machine learning with seven sparks from the causal revolution</a> . <i>CoRR</i> , abs/1801.04016.	
857		
858		
859	Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. 2016. <a href="#">Causal inference by using invariant prediction: identification and confidence intervals</a> . <i>Journal of the Royal Statistical Society: Series B (Statistical Methodology)</i> , 78(5):947–1012.	
860		
861		
862		
863		
864	Jonas Martin Peters, Dominik Janzing, and Bernhard Schölkopf. 2017. <i>Elements of Causal Inference: Foundations and Learning Algorithms</i> . MIT Press, Cambridge, MA, USA.	
865		
866		
867		
	Maxime Peyrard, Wei Zhao, Steffen Eger, and Robert West. 2021. <a href="#">Better than average: Paired evaluation of NLP systems</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2301–2315, Online. Association for Computational Linguistics.	868
		869
		870
		871
		872
		873
		874
		875
	Yusu Qian, Urwa Muaz, Ben Zhang, and Jae Won Hyun. 2019. <a href="#">Reducing gender bias in word-level language models with a gender-equalizing loss function</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop</i> , pages 223–228, Florence, Italy. Association for Computational Linguistics.	876
		877
		878
		879
		880
		881
		882
	Elan Rosenfeld, Pradeep Kumar Ravikumar, and Andrej Risteski. 2021. <a href="#">The risks of invariant risk minimization</a> . In <i>International Conference on Learning Representations</i> .	883
		884
		885
		886
	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. <a href="#">Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter</a> . <i>CoRR</i> , abs/1910.01108.	887
		888
		889
		890
	B. Schölkopf, F. Locatello, S. Bauer, N. R. Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio. 2021. <a href="#">Toward causal representation learning</a> . <i>Proceedings of the IEEE - Advances in Machine Learning and Deep Neural Networks</i> , 109(5):612–634.	891
		892
		893
		894
		895
	Bernhard Schölkopf. 2019. <a href="#">Causality for machine learning</a> . <i>CoRR</i> , abs/1911.10500.	896
		897
	Ilya Shpitser and Judea Pearl. 2008. <a href="#">Complete identification methods for the causal hierarchy</a> . <i>Journal of Machine Learning Research</i> , 9(64):1941–1979.	898
		899
		900
	Lifu Tu, Garima Lalwani, Spandana Gella, and He He. 2020. <a href="#">An empirical study on robustness to spurious correlations using pre-trained language models</a> . <i>Transactions of the Association for Computational Linguistics</i> , 8:621–633.	901
		902
		903
		904
		905
	Giorgos Vernikos, Katerina Margatina, Alexandra Chronopoulou, and Ion Androutsopoulos. 2020. <a href="#">Domain Adversarial Fine-Tuning as an Effective Regularizer</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3103–3112, Online. Association for Computational Linguistics.	906
		907
		908
		909
		910
		911
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923

- 924 Nanyang Ye, Kaican Li, Lanqing Hong, Haoyue Bai,  
925 Yiting Chen, Fengwei Zhou, and Zhenguo Li. 2021.  
926 [Ood-bench: Benchmarking and understanding out-of-](#)  
927 [distribution generalization datasets and algorithms.](#)
- 928 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Or-  
929 donez, and Kai-Wei Chang. 2017. [Men also like](#)  
930 [shopping: Reducing gender bias amplification using](#)  
931 [corpus-level constraints.](#) In *Proceedings of the 2017*  
932 *Conference on Empirical Methods in Natural Lan-*  
933 *guage Processing*, pages 2979–2989, Copenhagen,  
934 Denmark. Association for Computational Linguis-  
935 tics.
- 936 Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-  
937 Wei Chang. 2018. [Learning gender-neutral word em-](#)  
938 [beddings.](#) In *Proceedings of the 2018 Conference on*  
939 *Empirical Methods in Natural Language Processing*,  
940 pages 4847–4853, Brussels, Belgium. Association  
941 for Computational Linguistics.
- 942 Shanshan Zhao, Mingming Gong, Tongliang Liu, Huan  
943 Fu, and Dacheng Tao. 2020. [Domain generalization](#)  
944 [via entropy regularization.](#) In *Advances in Neural*  
945 *Information Processing Systems 33: Annual Confer-*  
946 *ence on Neural Information Processing Systems 2020,*  
947 *NeurIPS 2020, December 6-12, 2020, virtual.*
- 948 Kaiyang Zhou, Yongxin Yang, Timothy M. Hospedales,  
949 and Tao Xiang. 2020. [Learning to generate novel](#)  
950 [domains for domain generalization.](#) In *Computer*  
951 *Vision - ECCV 2020 - 16th European Conference,*  
952 *Glasgow, UK, August 23-28, 2020, Proceedings, Part*  
953 *XVI*, volume 12361 of *Lecture Notes in Computer*  
954 *Science*, pages 561–578. Springer.
- 955 Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang.  
956 2021. [Domain generalization with mixstyle.](#) In *9th*  
957 *International Conference on Learning Representa-*  
958 *tions, ICLR 2021, Virtual Event, Austria, May 3-7,*  
959 *2021.* OpenReview.net.



## A Illustration of iLM Architecture

In the main paper, we described formally the pseudo-code involved in training iLM models. The model architecture and the logic of the training schedule is illustrated in Fig. 4 for the special-case of 2 environments ( $n = 2$ ).

## B Details about Experiments

### B.1 Structured Noise Removal

**Data.** The data used for this experiment comes from an HTML Wikipedia Dump of August 2018. The files were pre-processed to remove the HTML content resulting in clean text articles. We randomly selected 6K articles with HTML (Env-B), and 6K different articles without HTML (Env-A). The test set contains 620 different articles without HTML.

**Hyper-parameters.** We ran the experiments reported in the main paper while varying several hyper-parameters: base transformers ( $\phi$ ): [distilBERT, ROBERTa], learning rates:  $[1e - 5, 5e - 5]$ , number of training steps: [10, 100, 200, 500, 2500, 5000], 5 random restarts with different random seeds,  $2 \cdot 2 \cdot 6 \cdot 5 = 120$ , ran with both eLM and iLM resulting in 240 experiments.

**Number of lines vs. number of articles.** In Fig. 2 of the main paper, we report the result of iLM and eLM when trained with environments having the same number of articles. However, the HTML articles have more lines and thus more *sentences*. Therefore, we also report in Fig. 5 the same analysis repeated when the number of lines between Env-A and Env-B is the same, meaning Env-B contains fewer articles. The conclusion remains largely unchanged in this scenario.

### B.2 Controlled Correlation Removal

**Data.** The dataset used for this experiment is Wikitext-2 (Merity et al., 2016) and the dictionary of gendered terms comes from Bordia and Bowman (2019) which was originally constructed to measure gender bias in language models.

The dictionary contains basic gender-pairs augmented with their variations in terms of casing, plural vs. singular forms and different spellings. The basic gendered pairs are: (actor, actress), (boy, girl), (boyfriend, girlfriend), (father, mother), (gentleman, lady), (grandson, granddaughter), (he, she),

	Unbalanced	Balanced
iLM ROBERTa	4.16	4.13
iLM distilBERT	5.82	5.81
eLM ROBERTa	4.14	4.14
eLM distilBERT	5.82	5.85

Table 1: Perplexities of iLM and eLM models after training.

(hero, heroine), (him, her), (husband, wife), (king, queen), (male, female), (man, woman), (mr., mrs.), (prince, princess), (son, daughter), (spokesman, spokeswoman), (stepfather, stepmother), (uncle, aunt)

**Hyper-parameters.** We ran the experiments reported in the main paper while varying several hyper-parameters: base-model ( $\phi$ ): [distilBERT, ROBERTa], learning-rates:  $[1e - 5, 5e - 5]$ , number of training steps: [10, 50, 100, 200, 1000, 2500], 5 random restarts with different random seeds  $2 \cdot 2 \cdot 6 \cdot 5 = 120$  experimental parameters, ran for both eLM and iLM for both the balanced and unbalanced setups resulting in 480 experiments.

**Details about the results.** Similar to the structured noise experiment, we report the performance of eLM and iLM as a function of the number of training steps and the probability that iLM is better than eLM when matched on hyper-parameter configuration. This is reported by Fig. 6 for two relative size: 25% (the modified environment has 4 times fewer examples) and 100%.

**Perplexities after training.** To ensure that the gender-based correlations were not removed at the cost of a worse perplexity, we report in Table 1 the perplexities of iLM models in comparison eLM ones on the test set of Wikitext-2. For reference, before our training distilBERT and ROBERTa had, this same test set, perplexities of 14.25 and 6.92, respectively.

In Table 1, the 95% confidence intervals all give uncertainties  $\approx 0.15$ , meaning that for a fixed base model (distilBERT or ROBERTa) all perplexities are within each other’s error bounds. There is no significant perplexity difference between eLM and iLM or between the unbalanced and balanced setups.

### B.3 Out-of-domain Generalization

**Data.** The data used for this experiment comes from subsamples of thePile (Gao et al., 2020). Af-

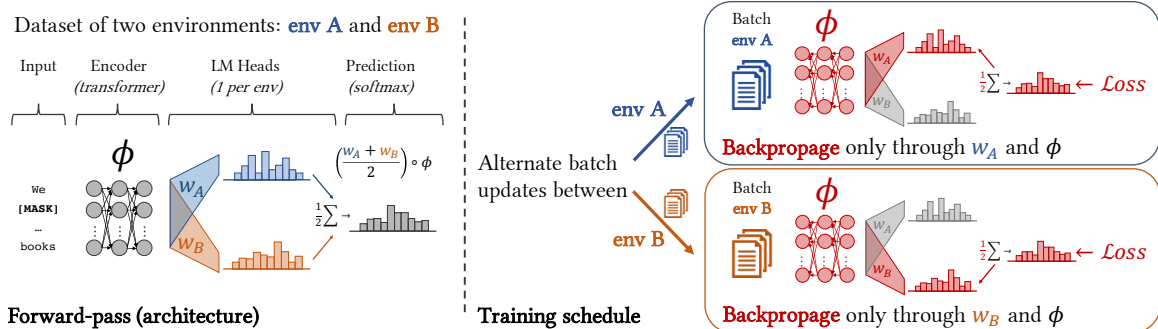


Figure 4: **Model description** In the forward pass, input text goes through the main body of language model noted  $\phi$  (e.g., a Transformer (Devlin et al., 2019)), then one head per environment predicts logits over the vocabulary. These predictions are averaged over all heads and go through a Softmax. During training, the model receives a batch of data from one environment  $e$  and performs a gradient update only on the parameters of the main body of the language model ( $\phi$ ) and on the parameters of the head tied to this environment  $w_e$ . Then batches are taken from each environment in a round-robin fashion.

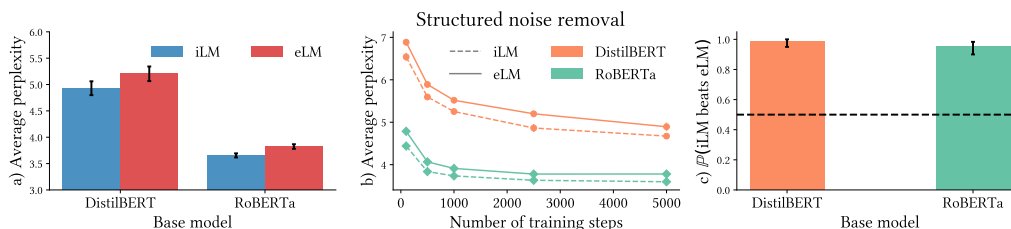


Figure 5: Structured noise removal experiment with environments having the same number of lines: a) average perplexity over all hyper-parameters b) average perplexity as a function of the number of training steps (for learning rate  $10^{-5}$ ), c) Probability that iLM is better than eLM when compared on the same hyper-parameters

1046 ter the result of our sampling described in the main  
1047 paper, 8 domains have ended-up as test domain.

1048 **Hyper-parameters.** We ran the experiments re-  
1049 ported in the main paper while varying several  
1050 hyper-parameters: base-model ( $\phi$ ): [distilBERT,  
1051 RoBERTa], learning-rates: [ $1e-5, 5e-5$ ], number  
1052 of training steps: [100, 1000, 2500, 5000], number  
1053 of environments for training: [3, 9, 13], 5 random  
1054 restarts with different random seeds and different  
1055 choices of training/testing domains.

1056 In Fig. 7, we report the probability that iLM  
1057 has lower perplexity than eLM as a function of  
1058 the number of training steps in Fig. 7 (a) and as a  
1059 function of the number of training environments  
1060 Fig. 7 (b).

1061 We observe that overall iLM is better perplexi-  
1062 ties on unseen domains. The advantage of iLM in-  
1063 creases with the number of training steps (Fig. 7 a)  
1064 but also with number of training environments  
1065 (Fig. 7 b). This indicates that using more envi-  
1066 ronments is even more beneficial for iLM than for  
1067 eLM.

	iLM	eLM
arxiv	<b>5.71</b>	5.93
openwebtext	3.90	3.96
pile-cc	4.42	4.44
uspto	<b>4.14</b>	4.19
pubmed-abstract	4.13	4.17
pubmed-central	<b>4.23</b>	4.29
github	<b>5.84</b>	5.93
youtube	4.78	4.76

Table 2: Perplexities of iLM and eLM models for both RoBERTa on testing domains subsampled from thePile. The bold font indicates that iLM is significantly better than eLM ( $p < .05$  paired t-test).

1068 **Perplexities.** In the main paper, we focus on the  
1069 paired comparison between iLM and eLM. In Ta-  
1070 ble 2, we report the test perplexities of iLM and  
1071 eLM for distilBERT and RoBERTa average over  
1072 different hyper-parameters. We observe that differ-  
1073 ences between eLM and iLM are smaller than for  
1074 other experiments but iLM still has advantage over  
1075 eLM.

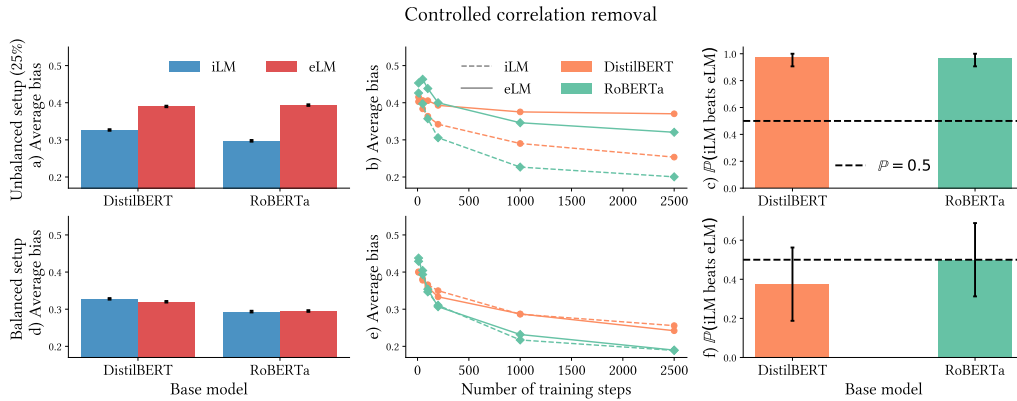


Figure 6: **Controlled correlation removal experiment:** On the first row, the modified environment is 25% of the size of the unmodified environment. On the second row, both have the same number of samples. On the left-most column, average bias over all hyper-parameters. On the center column: average bias as a function of the number of training steps. On the right-most column: Probability that iLM is less biased than eLM when compared on the same hyper-parameters.

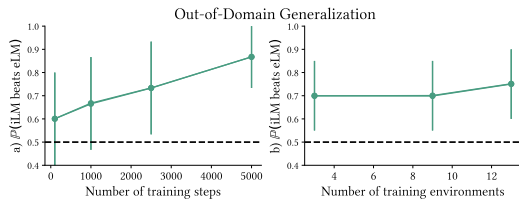


Figure 7: **OoD generalization:** a) Probability that iLM is better than eLM all hyper-parameters being the same as a function of: the number of training steps in a) and the number of training environments in b).

#### B.4 Languages as Environments

One question that might arise from iLM training schedule is whether it simply focuses on surface-level lexical correlations in the data. For example, if the lexical correlations are different across environments, maybe no correlation remain generalizable and iLM learns an empty set of correlations. To better demonstrate that iLM operate on latent variable and not on surface-level correlations, we perform a simple experiment with languages as environments.

**Description.** We use two languages with no lexical overlap: English and Farsi. We put english Wikipedia articles as one environment and farsi Wikipedia articles as the other. In this setup, no surface-level correlations can generalize across environment as the two environments don't even have the same vocabulary.

We train iLM with a multilingual pre-trained RoBERTa: XLM-RoBERTa for 5000 steps with these two environments of equal size (10K arti-

cles per language). Then, we test whether this choice of environments destructs previously learn correlations in the language model by comparing perplexities on a balanced held-out test set of english and farsi documents against the model before finetuning. If the perplexities decrease, we would conclude that iLM destroy surface-level correlations.

**Results.** We found that before finetuning, XLM-RoBERTa had a perplexity of 14.56 on the held-out test set, where iLM could improve it perplexity down to 6.44. This indicates that iLM with environments having no lexical overlap does not destroy previously learned correlations. It can even improve its perplexities for each language. A possible reason why iLM can even improve so dramatically compared to before finetuning might come from the fact that  $\phi$  learns to recognize the languages, separate them and treat them separately. Similar effects have been observed in previous work (Guo et al., 2021) when the correlation between the environment index and the target variable is very strong (which is the case here).

#### B.5 Head dynamics

The main components of our framework are the heads and their training dynamic. Therefore, we investigate aspects related to behaviour of the heads.

**Description.** During training, the loss of each head is still entangled with the prediction of every other head. So we wonder whether the heads still capture information related to the environment it is tied to during training. In particular, we ask (i) whether the

1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128



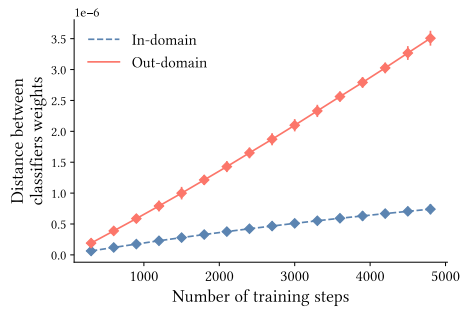


Figure 8: Comparing distance between heads weights in- and out-domain as functions of the number of training step. (95% confidence interval from random restart with different seeds.)

parameters of the heads for different environments are drifting apart during training? Indeed, all heads are initialized to the same pretrained weights at the beginning of training. (ii) Are the parameters of the heads predicting which environments are more similar?

**Experimental setup.** To answer these two questions in one go, we take two environments  $A$  and  $B$  and split each of them into two new environments resulting in  $A_1, A_2, B_1,$  and  $B_2$  such that  $A_1$  and  $A_2$  are very similar  $B_1$  and  $B_2$  are very similar but  $A_i$  and  $B_i$  are different. We then train iLM with the four environments and, thus, with four heads  $w_{A_1}, w_{A_2}, w_{B_1},$  and  $w_{B_2}$ . We measure whether the heads’ weights can predict the similarities between  $A$ ’s and  $B$ ’s environments.

$$D_{in} = \frac{1}{2} (d(w_{A_1}, w_{A_2}) + d(w_{B_1}, w_{B_2})), \quad (3)$$

$$D_{out} = \frac{1}{4} \sum_{i,j} d(w_{A_i}, w_{B_j}), \quad (4)$$

where  $d$  is the L2 distance between the linearized weights of two heads. Then,  $D_{in}$  is the average distance between heads tied the same domain, and  $D_{out}$  is the average distance between heads tied to different domains. Remember that in this case, there are 2 domains  $A$  and  $B$  and 4 environments  $A_i$  and  $B_j$ .

In this experiment, we randomly select the base environments  $A$  and  $B$  from the domains of thePile ( $A$  is the Enron-Email, and  $B$  is PubMed abstract). We create  $A_i$  and  $B_j$  by randomly subsampling 2 environments of the same size from each domain. We train iLM with ROBERTa for 5000 training steps, taking checkpoints of the heads every 500 steps. We perform 10 random restarts with different seeds to uncertainty estimates. In Fig. 8, we report  $D_{in}$

and  $D_{out}$  as functions of the number of training steps.

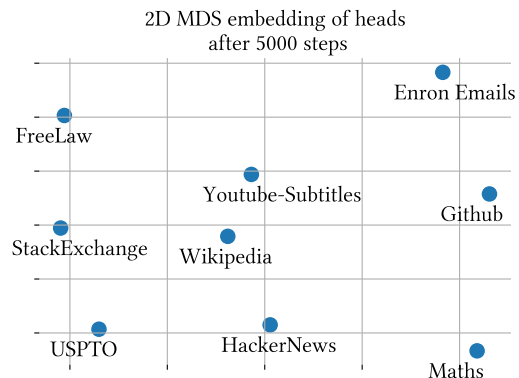


Figure 9: Heads embeddings: 2D projection of the heads parameters similarity structure after training iLM with ROBERTa for 5000 steps with 9 domains. Each dot represent one head of the model after training and the labels indicate to which domain it is tied to.

**Analysis.** We first notice that indeed the heads are drifting apart from each other as training advances. More interestingly, the distance between heads from the same domain is significantly much smaller than the distance between heads from different domains. We conclude that heads retain environment-specific information in their parameters and are predictive of environment similarities.

Now, we visualize the geometry of head similarity by training iLM with ROBERTa for 5000 steps with 9 environments from thePile: . After training, we take the heads’ parameters and compute the pairwise distance between all 9 heads and embed them in 2D with Multi-Dimensional Scaling to visualize the similarity structure. The result is depicted in Fig. 9.