Are Larger Pretrained Language Models Uniformly Better? Comparing Performance at the Instance Level

Ruiqi Zhong Dhruba Ghosh Dan Klein Jacob Steinhardt

Computer Science Division, University of California, Berkeley {ruiqi-zhong, djghosh13, klein, jsteinhardt}@berkeley.edu

Abstract

Larger language models have higher accuracy on average, but are they better on every single instance (datapoint)? Some work suggests larger models have higher out-ofdistribution robustness, while other work suggests they have lower accuracy on rare subgroups. To understand these differences, we investigate these models at the level of individual instances. However, one major challenge is that individual predictions are highly sensitive to noise in the randomness in training. We develop statistically rigorous methods to address this, and after accounting for pretraining and finetuning noise, we find that our BERT-LARGE is worse than BERT-MINI on at least 1-4% of instances across MNLI, SST-2, and QQP, compared to the overall accuracy improvement of 2-10%. We also find that finetuning noise increases with model size, and that instance-level accuracy has momentum: improvement from BERT-MINI to BERT-MEDIUM correlates with improvement from BERT-MEDIUM to BERT-LARGE. Our findings suggest that instance-level predictions provide a rich source of information; we therefore recommend that researchers supplement model weights with model predictions.

1 Introduction

Historically, large deep learning models (Peters et al., 2018; Devlin et al., 2019; Lewis et al., 2020; Raffel et al., 2019) have improved the state of the art on a wide range of tasks and leaderboards (Schwartz et al., 2014; Rajpurkar et al., 2016; Wang et al., 2018), and empirical scaling laws predict that larger models will continue to increase performance (Kaplan et al., 2020). However, little is understood about such improvement at the instance (datapoint) level. Are larger models uniformly better? In other words, are larger pretrained models better at every instance, or are they better at some instances, but worse at others?

Prior works hint at differing answers. Hendrycks et al. (2020) and Desai and Durrett (2020) find that larger pretrained models consistently improve out-of-distribution performance, which implies that they might be uniformly better at a finer level. Henighan et al. (2020) claim that larger pretrained image models have lower downstream classification loss for the majority of instances, and they predict this trend to be true for other data modalities (e.g. text). On the other hand, Sagawa et al. (2020) find that larger non-pretrained models perform worse on rare subgroups; if this result generalizes to pretrained language models, larger models will not be uniformly better. Despite all the indirect evidence, it is still inconclusive how many instances larger pretrained models perform worse

A naïve solution is to finetune a larger model, compare it to a smaller one, and find instances where the larger model is worse. However, this approach is flawed, since model predictions are **noisy at the instance level**. On MNLI in-domain development set, even the same architecture with different finetuning seeds leads to different predictions on ~8% of the instances. This is due to under-specification (D'Amour et al., 2020), where there are multiple different solutions that can minimize the training loss. Since the accuracy improvement from our BERT-BASE¹ to BERT-LARGE is 2%, most signals across different model sizes will be dominated by noise due to random seeds.

To account for the noise in pretraining and finetuning, we define *instance accuracy* as "how often a model correctly predicts an instance" (Figure 1 left) in expectation across pretraining and finetuning seeds. We estimate this quantity by pretraining 10 models with different seeds, finetuning 5 times for each pretrained models (Figure 1 middle), and

¹This is not the original release by Devlin et al. (2019); we pretrained models ourselves.

averaging across them.

However, this estimate is still inexact, and we might falsely observe smaller models to be better at some instances by chance. Hence, we propose a random baseline to estimate the fraction of *false discoveries* (Section 3, Figure 1 right) and formally upper-bound the false discoveries in Section 4. Our method provides a better upper bound than the classical Benjamini-Hochberg procedure with Fisher's exact test.

Using the 50 models for each size and our improved statistical tool, we find that, on the MNLI in-domain development set, the accuracy "decays" from BERT-LARGE to BERT-MINI on at least ~4% of the instances, which is significant given that the improvement in overall accuracy is 10%. These decaying instances contain more controversial or wrong labels, but also correct ones (Section 4.2). Therefore, larger pretrained language models are not uniformly better.

We make other interesting discoveries at the instance level. Section 5 finds that instance-level accuracy has momentum: improvement from MINI to MEDIUM correlates with improvement from MEDIUM to LARGE. Additionally, Section 6 attributes variance of model predictions to pretraining and finetuning random seeds, and finds that finetuning seeds cause more variance for larger models. Our findings suggest that instance-level predictions provide a rich source of information; we therefore recommend that researchers supplement model weights with model predictions. In this spirit, we release all the pretrained models, model predictions, and code here: https://github.com/ruigi-zhong/ac12021-instance-level.

2 Data, Models, and Predictions

To investigate model behavior, we considered different sizes of the BERT architecture and finetuned them on Quora Question Pairs (QQP²), Multi-Genre Natural Language Inference (MNLI; Williams et al. (2020)), and the Stanford Sentiment Treebank (SST-2; Socher et al. (2013)). To account for pretraining and finetuning noise, we averaged over multiple random initializations and training data order, and thus needed to pretrain our own models rather than downloading off the internet. Following Turc et al. (2019) we trained 5 architectures of increasing size: MINI

(L4/H256, 4 Layers with hidden dimension 256), SMALL (L4/H512), MEDIUM (L8/H512), BASE (L12/H768), and LARGE (L24/H1024). For each architecture we pre-trained models with 10 different random seeds and fine-tuned each of them 5 times (50 total) on each task; see Figure 1 middle. Since pretraining is computationally expensive, we reduced the context size during pretraining from 512 to 128 and compensated by increasing training steps from 1M to 2M. Appendix A includes more details about pretraining and finetuning and their computational cost, and Appendix B verifies that our cost-saving changes do not affect accuracy qualitatively.

Notation. We use i to index an instance in the evaluation set, s for model sizes, P for pretraining seeds and F for finetuning seeds. c is a random variable of value 0 or 1 to indicate whether the prediction is correct. Given the pretraining seed P and the finetuning seed F, $c_i^s = 1$ if the model of size s is correct on instance i, 0 otherwise. To keep the notation uncluttered, we sometimes omit these superscripts or subscripts if they can be inferred from context.

Unless otherwise noted, we present results on the MNLI in-domain development set in the main paper.

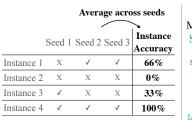
3 Comparing Instance Accuracy

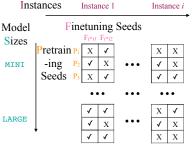
To find the instances where larger models are worse, a naïve approach is to finetune a larger pretrained model, compare it to a smaller one, and find instances where the larger is incorrect but the smaller is correct. Under this approach, BERT-LARGE is worse than BERT-BASE on 4.5% of the instances and better on 7%, giving an overall accuracy improvement of 2.5%.

However, this result is misleading: even if we compare two BERT-BASE model with different finetuning seeds, their predictions differ on 8% of the instances, while their accuracies differ only by 0.1%; Table 1 reports this baseline randomness across model sizes. Changing the pretraining seed also changes around 2% additional predictions beyond finetuning.

Table 1 also reports the standard deviation of overall accuracy, which is about 40 times smaller. Such stability starkly contrasts with the noisiness at the instance level, which poses a unique challenge.

²https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs





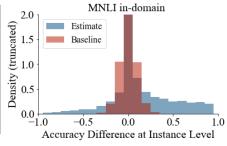


Figure 1: **Left**: Each column represents the same architecture trained with a different seed. We calculate accuracy for each instance (row) by averaging across seeds (column), while it is usually calculated for each model by averaging across instances. **Middle**: A visual layout of the model predictions we obtain, which is a binary-valued tensor with 4 axes: model size s, instance i, pretraining seeds P and finetuning seeds F. **Right**: for each instance, we calculate the accuracy gain from MINI to LARGE and plot the histogram in blue, along with a random baseline in red. Since the blue distribution has a bigger left tail, smaller models are better at some instances.

	$\mathrm{Diff}_{\mathrm{FTune}}$	$\mathrm{Diff}_{\mathrm{PTrain}}$	Std_{all}
MINI	7.2%	10.7%	0.2%
SMALL	7.2%	10.7%	0.3%
MEDIUM	8.0%	10.7%	0.3%
BASE	8.5%	10.6%	0.2%
LARGE	8.6%	10.1%	0.2%

Table 1: Larger model sizes are at the bottom rows. Diff $_{\rm FTune}$: how much do the predictions differ, if two models have the same pretraining seed but different finetuning seeds F? Diff $_{\rm PTrain}$: the difference if the pretraining seeds P are different. Std $_{\rm all}$: the standard deviation of overall accuracy, around 40 times smaller than Diff $_{\rm FTune}$.

Instance-Level Metrics To reflect this noisiness, we define the *instance accuracy* Acc_i^s to be how often models of size s predict instance i correctly,

$$Acc_i^s := \mathbb{E}_{P,F}[c_i^s]. \tag{1}$$

The expectation is taken with respect to the pretraining and finetuning randomness P and F. We estimate Acc_i^s via the empirical average Acc_i^s accross 10 pretraining \times 5 finetuning runs.

We histogram $\hat{\mathrm{Acc}}_i^s$ in Figure 2 (a). On most instances the model always predicts correctly or incorrectly ($\hat{\mathrm{Acc}}=0$ or 1), but a sizable fraction of accuracies lie between the two extremes.

Recall that our goal is to find instances where larger models are less accurate, which we refer to as *decaying* instances. We therefore study the *instance difference* between two model sizes s_1 and s_2 , defined as

$$_{s_2}^{s_1} \Delta \operatorname{Acc}_i := \operatorname{Acc}_i^{s_2} - \operatorname{Acc}_i^{s_1}, \tag{2}$$

which is estimated by the difference between the

accuracy estimates \hat{Acc}_i^s , i.e.

$$_{s_2}^{s_1} \Delta \hat{\operatorname{Acc}}_i := \hat{\operatorname{Acc}}_i^{s_2} - \hat{\operatorname{Acc}}_i^{s_1}. \tag{3}$$

We histogram $_{\text{LARGE}}^{\text{BASE}}\Delta\hat{A}cc_{i}$ in Figure 2 (b). We observe a unimodal distribution centered near 0, with tails on both sides. Therefore, the estimated differences for some instances are negative.

However, due to estimation noise, we might falsely observe this accuracy decay by chance. Therefore, we introduce a random baseline ${}^{s_1}_{s_2}\Delta Acc'$ to control for these false discoveries. Recall that we have 10 smaller pretrained models and 10 larger ones. Our baseline splits these into a group A of 5 smaller + 5 larger, and another group B of the remaining 5+5. Then the empirical accuracies \hat{Acc}^A and \hat{Acc}^B are identically distributed, so we take our baseline ${}^{s_1}_{s_2}\Delta Acc'$ to be the difference $\hat{Acc}^A - \hat{Acc}^B$. We visualize and compare how to calculate ${}^{s_1}_{s_2}\Delta \hat{Acc}$ and ${}^{s_1}_{s_2}\Delta Acc'$ in Figure 3.

We histogram this baseline $^{\text{BASE}}_{\text{LARGE}}\Delta\text{Acc}'$ in Figure 2 (b), and find that our noisy estimate $^{\text{BASE}}_{\text{LARGE}}\Delta\hat{\text{Acc}}$ has a larger left tail than the baseline. This suggests that decaying instances exist. We similarly compare MINI to LARGE in Figure 2 (c) and find an even larger left tail.

4 Quantifying the Decaying Instances

The left tail of $\Delta \hat{A}cc$ noisily estimates the fraction of decaying instances, and the left tail of the random baseline $\Delta Acc'$ counts the false discovery fraction due to the noise. Intuitively, the true fraction of decaying instances can be captured by the difference of these left tails, and we formally quantify this below.

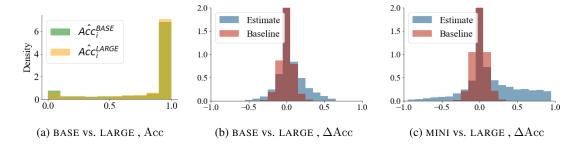


Figure 2: (a) The distribution of instance accuracy \hat{Acc}_i . (b, c) Histogram of instance difference estimate (x-axis), \hat{Acc}_i (blue) and its baseline \hat{AAcc}_i (rad) compares BASE and LABGE. To better visualize, we truncated the eare indeed instances

			Â	cc ^{LA}	RGE	$= 0.75 - \hat{Acc}^{A}$	=	$_{LARGE}^{MINI}\Delta\hat{A}cc=.33$			
LAF	RGE		F(*)1	F(*)	F _{(*)3}	MIN	Ι	F(*)1	F(*)	F(*)3	
D	I	1	Х	√	√	Group A	P_1	Х	√	√	$\hat{Acc}^A = 0.58$
Pretrain-	I	2	√	√	Х	Group A	P ₂	√	Х	Χ	Acc = 0.58
ing Seeds	I	3	√	√	Х	Crown D	P ₃	Х	Х	Х	$\hat{Acc}^B = 0.58$
Secus	I	4	√	√	√	Group B	P ₄	Χ	√	√	Acc = 0.58
											$_{Large}^{Min i} \Delta Acc' = 0$

Figure 3: The tables are model predictions with visual notations established in Figure 1 middle. $\Delta \hat{A}cc$ (blue) is the mean difference between the left and the right table, each corresponding to a model size. The random baseline $\Delta Acc'$ (red) is the mean difference between group A (orange) cells and group B (green), which are identically and independently distributed.

Suppose instance i is drawn from the empirical evaluation distribution. Then we can define the true decaying fraction Decay as

$$Decay := \mathbb{P}_i[\Delta Acc_i < 0]. \tag{4}$$

Since $\Delta \mathrm{Acc}_i$ is not directly observable and $\Delta \mathrm{Acc}_i$ is noisy, we add a buffer and only consider instances with $\Delta \mathrm{Acc}_i \leq t$, which makes it more likely (but still uncertain) that the true $\Delta \mathrm{Acc}_i < 0$. We denote this "discovery fraction" $\mathrm{Decay}(t)$ as

$$\hat{\text{Decay}}(t) := \mathbb{P}_i[\Delta \hat{\text{Acc}}_i \le t].$$
 (5)

Similarly, we define a baseline control (false discovery fraction) $\operatorname{Decay}'(t) := \mathbb{P}_i[\Delta \operatorname{Acc}_i' \leq t]$. Hence, Decay and Decay' are the cumulative distribution function of $\Delta \operatorname{Acc}$ and $\Delta \operatorname{Acc}'$ (Figure 4).

We have the following theorem, which we formally state and prove in Appendix D:

Theorem 1 (Informal) If all the random seeds are independent, then for all thresholds t,

$$\operatorname{Decay} \ge \mathbb{E}[\hat{\operatorname{Decay}}(t) - \operatorname{Decay}'(t)]$$
 (6)

Proof Sketch Suppose we observe $c_{R_{1...2k}}^{s_1}$ and $c_{R_{2k+1...4k}}^{s_2}$, where there are 2k different random seeds for each model size 3 . Then

$$\Delta \hat{A} cc_i := \frac{1}{2k} \left(\sum_{j=1}^{2k} c_{R_j,i}^{s_1} - \sum_{j=2k+1}^{4k} c_{R_j,i}^{s_2} \right), \quad (7)$$

and hence the discovery rate $\hat{\mathrm{Decay}}(t)$ is defined as

$$\hat{\text{Decay}}(t) := \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} \mathbf{1}[\Delta \hat{\text{Acc}} \le t].$$
 (8)

For the random baseline estimator, we have

$$\Delta \operatorname{Acc}_{i}' := \frac{1}{2k} \left(\sum_{j=1}^{k} c_{R_{j},i}^{s_{1}} + \sum_{j=2k+1}^{3k} c_{R_{j},i}^{s_{2}} - \sum_{j=k+1}^{2k} c_{R_{j},i}^{s_{1}} - \sum_{j=3k+1}^{4k} c_{R_{j},i}^{s_{2}} \right),$$
(9)

and the false discovery control Decay' is defined as

$$Decay'(t) := \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} \mathbf{1}[\Delta Acc_i' \le t].$$
 (10)

Formally, the theorem states that

$$\operatorname{Decay} \geq \mathbb{E}_{R_1...R_{4k}}[\hat{\operatorname{Decay}}(t) - \operatorname{Decay}'(t)],$$
 (11) which is equivalent to

$$\sum_{i=1}^{|\mathcal{T}|} (\mathbf{1}[\Delta Acc_i < 0] - \mathbb{P}[\Delta \hat{A}cc_i \le t] + \mathbb{P}[\Delta Acc_i' \le t]) \ge 0$$
(12)

³We assumed even number of random seeds since we will mix half of the models from each size to compute the random baseline

Cumulative Distribution of ΔAcc , Mini vs. Large

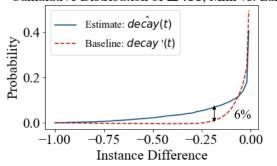


Figure 4: The cumulative distribution function of the histogram in Figure 2 (c); only the negative x-axis is shown because it corresponds to decays. The maximum difference between the two curves (6%) is a lower bound of the true decaying fraction.

Hence, we can declare victory if we can prove that for all i, if $\Delta Acc_i \ge 0$,

$$\mathbb{P}[\Delta Acc_i' \le t] \ge \mathbb{P}[\Delta \hat{A}cc_i \le t].$$

This is easy to see, since $\Delta \mathrm{Acc}_i'$ and $\Delta \mathrm{\hat{A}cc}_i$ are both binomial distributions with the same n, but the first has a larger rate. 4

Roughly speaking, the true decaying fraction is at least the difference between $\hat{\text{Decay}}(t)$ and $\hat{\text{Decay}}'(t)$ at every threshold t. Therefore, we take the maximum difference between $\hat{\text{Decay}}(t)$ and $\hat{\text{Decay}}'(t)$ to lower-bound the fraction of decaying instances. For example, Figure 4 estimates the true decaying fraction between MINI and LARGE to be at least 6%.

We compute this lower bound for other pairs of model sizes in Table 2, and the full results across other tasks and model size pairs are in Appendix C. In all of these settings we find a non-zero fraction of decaying instances, and larger model size differences usually lead to more decaying instances.

Unfortunately, applying Theorem 1 as above is not fully rigorous, since some finetuning runs share the same pretraining seeds and hence are dependent.⁶ To obtain a statistically rigorous lower bound, we slightly modify our target of interest. Instead of examining individual finetuning runs, we ensemble our model across 5 different finetuning runs for each pretraining seed; these predictions

$s_1 \setminus s_2$	MINI	SMALL	BASE	LARGE
91 / 92	IVIIIVI	SMALL	DASE	LAKUL
MINI	N/A	9%	18%	21%
SMALL	3%	N/A	14%	18%
BASE	6%	5%	N/A	10%
LARGE	6%	5%	2%	N/A

Table 2: We lower-bound the fraction of instances that improve when model size changes from s_1 (row) to s_2 (column). For example, when model size decreases from LARGE to MINI, 6% of instances improve (i.e. decays).

Threshold	Decay	Decay'	Diff
t = -0.4	4.22%	3.49e - 3	3.87%
• • •			
t = -0.9	0.91%	$1.44e{-7}$	0.91%
t = -1.0	0.48%	$2.06e{-8}$	0.48%

Table 3: Comparing MINI vs. LARGE by calculating the discovery fraction $\hat{\text{Decay}}$, the false discovery control $\hat{\text{Decay}}$, and their difference (Diff) under different thresholds t. LARGE is worse on at least $\sim 4\%$ (maximum Diff) of instances.

are essentially the same as individual finetuning runs, except that the finetuning randomness is averaged out. Hence we obtain 10 independent sets of model predictions with different random seeds, which allows us to apply Theorem 1.

We compare MINI to LARGE using these ensembles and report the discovery Decay and the baseline Decay' in Table 3. Taking the maximum difference across thresholds, we estimate at least $\sim 4\%$ of decaying instances. This estimate is lower than the previous 6% estimate, which used the full set of 50 models' predictions assuming they were independent. However, this is still a meaningful amount, given that the overall accuracy improvement from MINI to LARGE is 10%.

4.1 Fisher's Test + Benjamini-Hochberg

Here is a more classical approach to lower-bound the decaying fraction. For each instance, we compute a significance level α under the null hypothesis that the larger model is better, using Fisher's exact test. We sort the significance levels ascendingly, and call the p^{th} percentile α_p . Then we pick a false discovery rate q (say, 25%), find the largest p s.t. $\alpha_p < pq$, and estimate the decaying fraction to be at least p(1-q). This calculation is known as the Benjamini-Hochberg procedure (Benjamini and Hochberg, 1995).

To compare our method with this classical ap-

⁴More details are in Appendix D.

 $^{^5}$ Adaptively picking the best threshold t depending on the data may incur a slight upward bias. Appendix E estimates that the relative bias is at most 10% using a bootstrap method.

⁶Although we anticipate such dependencies do not cause a substantial difference, as discussed in Appendix D.1.

s1	s2		2	6	10
MINI	LARGE	ours	1.9%	3.1%	4.0%
MINI	LARGE	BH	0.0%	0.9%	1.9%
BASE	LARGE	ours	0.4%	0.9%	1.2%
BASE	LARGE	BH	0.0%	0.0%	0.0%

Table 4: We compare our method to the Fisher's exact test + Benjamin-Hochberg (BH) procedure described in Section 4. For all different model size pairs and number of pretrained models available, ours always provides a higher (better) lower bound of the decaying fraction.

proach, we estimate the lower bound of the decaying fraction for different pairs of model sizes with different numbers of pretrained models available. To make sure our choice of the false discovery rate q does not bias against the classical approach, we adaptively choose q to maximize its performance. Appendix F includes the full results and Table 4 is a representative subset.

We find that our approach is more powerful, particularly when the true decaying fraction is likely to be small and only a few models are available, which is usually the regime of interest. For example, across all pairs of model sizes, our approach only needs 2 random seeds (i.e. pretrained models) to provide a non-zero lower bound on the decaying fraction, while the classical approach sometimes fails to do this even with 10 seeds. Intuitively, when fewer seeds are available, the smallest possible significance level for each instance is larger than the decaying fraction, hence hurting the classical approach.

4.2 Understanding the Decaying Instances

We next manually examine the decaying instances to see whether we can find any interpretable patterns. One hypothesis is that all the decaying fractions are in fact mislabeled, and hence larger models are not in fact worse on any instances.

To investigate this hypothesis, we examined the group of instances where $^{\text{MINI}}_{\text{LARGE}}\Delta\hat{A}\text{cc}_i \leq -0.9$. MINI is almost always correct on these instances, while LARGE is almost always wrong, and the false discovery fraction is tiny. For each instance, we manually categorize it as either: 1) Correct, if the label is correct, 2) Fine, if the label might be controversial but we could see a reason why this label is reasonable, 3) Wrong, if the label is wrong, or 4) Unsure, if we are unsure about how to label this instance. Each time we annotate, with 50% probability we randomly sample either a decaying

	Correct	Fine	Wrong	Unsure
$MNLI^D$	66%	17%	9%	5%
$MNLI^C$	86%	5%	5%	1%
SST- 2^D	55%	8%	10%	25%
SST- 2^C	88%	4%	0%	6%
QQP^D	60%	26%	10%	2%
QQP^C	87%	10%	1%	0%

Table 5: MINI vs. LARGE . We examine whether there are mislabels for the **D**ecaying fractions (superscript D) and the rest of the dataset (Control group C). The decaying fraction contains more mislabels, but includes correct labels as well.

instance or an instance from the remaining dataset as a control. We are blind to which group it comes from

For each task of MNLI, QQP, and SST-2, the first author annotated 100 instances (decay + control group) (Table 5). We present all the annotated decaying instances in Appendix J.

Conclusion We find that the decaying fraction has more wrong or controversial labels, compared to the remaining instances. However, even after we adjust for the fraction of incorrect labels, the Decay fraction still exceeds the false discovery control. This implies that MINI models are better than LARGE models on some correctly labeled instances. The second author followed the same procedure and reproduced the same qualitative results.

However, we cannot find an interpretable pattern for these correctly labeled decaying instances by simply eyeballing. We discuss future directions to discover interpretable categories in Section 7.

5 Correlation of Instance Difference

We next investigate whether there is a momentum of instance accuracy increase: for example, if the instance accuracy improves from $\operatorname{MINI}(s_1)$ to $\operatorname{MEDIUM}(s_2)$, is it more likely to improve from $\operatorname{MEDIUM}(s_2)$ to $\operatorname{LARGE}(s_3)$?

The naïve approach is to calculate the Pearson correlation coefficient between $^{\rm MINI}_{\rm MEDIUM}\Delta\hat{A}cc$ and $^{\rm MEDIUM}_{\rm LARGE}\Delta\hat{A}cc$, and we find the correlation to be zero. However, this is partly an artifact of accuracies being bounded in [0,1]. If MEDIUM drastically improves over MINI from 0 to 1, there is no room for LARGE to improve over MEDIUM. To remove this inherent negative correlation, we calculate the correlation conditioned on the accuracy of the middle-

$(s_1, s_2, s_3) \downarrow \text{Buckets} \rightarrow$	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	1.00
SMALL, MEDIUM, BASE										
MINI, MEDIUM, LARGE	0.03	0.15	0.18	0.33	0.17	0.16	0.22	0.20	0.19	0.09

Table 6: Each row corresponds to a triplet of model sizes. Each column t represents a bucket that contains instances with $\hat{\mathrm{Acc}}^{s_2} \in [t-0.1,t]$. Within each bucket, we calculate the Pearson correlation coefficient between the estimated accuracy improvements: $\frac{s_1}{s_2}\Delta\hat{\mathrm{Acc}}$ and $\frac{s_2}{s_3}\Delta\hat{\mathrm{Acc}}$. These correlations are positive and become higher when model size differences are small.

sized model, $\hat{\text{Acc}}^{\text{MEDIUM}}$.

Therefore, we bucket instances by their estimated MEDIUM accuracy into intervals of size 0.1, and we find the correlation to be positive within each bucket (Table 6, row 2). This fixes the problem with the naïve approach by getting rid of the negative correlation, which could have misled us to believe that improvements by larger models are uncorrelated.

We additionally find that the correlations between improvements become stronger when model size differences are smaller. Table 6 row 1 reports results for another model size triplet with smaller size difference, i.e. $(s_1, s_2, s_3) = (SMALL, MEDIUM, BASE)$, and the correlation is larger for all buckets. Results for more tasks and size triplets are in Appendix G and the same conclusions hold qualitatively.

6 Variance at the Instance Level

Section 3 found that the overall accuracy has relatively low variance, but model predictions are noisy. This section formally analyzes variance at the instance level. For each instance, we decompose its loss into three components: ${\rm Bias}^2$, variance due to pretraining randomness, and variance due to finetuning randomness. Formally, we consider the 0/1 loss:

$$\mathcal{L}_i := 1 - c_i = (1 - c_i)^2,$$
 (13)

where c_i is a random variable 0/1 indicating whether the prediction is correct or incorrect, with respect to randomness in pretraining and finetuning. Therefore, by bias-variance decomposition and total variance decomposition, we have

$$\mathcal{L}_i = \operatorname{Bias}^2_i + \operatorname{PretVar}_i + \operatorname{FineVar}_i, \quad (14)$$

where, by using P and F as pretraining and finetuning random seeds:

$$\operatorname{Bias}^{2}_{i} := (1 - \mathbb{E}_{P,F}[c_{i}])^{2}, \qquad (15)$$

$$\operatorname{PretVar}_{i} := \operatorname{Var}_{P}[\mathbb{E}_{F}[c_{i}]],$$

$$\operatorname{FineVar}_{i} := \mathbb{E}_{P}[\operatorname{Var}_{F}[c_{i}]],$$

	Bias^2	PretVar	FineVar
MINI	0.203	0.017	0.036
SMALL	0.179	0.017	0.036
MEDIUM	0.157	0.014	0.040
BASE	0.134	0.010	0.043
LARGE	0.111	0.007	0.043

Table 7: The bias, pretraining variance, and finetuning variance for each model size, averaged across all test instances. Finetuning variance is much larger than pretraining variance; larger models have larger finetuning variance.

capturing "how wrong is the average prediction", variance due to pretraining, and variance due to finetuning seeds, respectively.

We can directly estimate FineVar by first calculating the sample variance across finetuning runs for each pretraining seed, and then averaging the variances across the pretraining seeds. Estimating PretVar is more complicated. A naïve approach is to calculate the empirical variance, across pretraining seeds, of the average accuracy across finetuning seeds. However, the estimated average accuracy for each pretraining seed is noisy itself, which causes an upward bias on the PretVar estimate. We correct this bias by estimating the variance of the estimated average accuracy and subtracting it from the naïve estimate; see Appendix H for details, as well as a generalization to more than two sources of randomness. Finally, we estimate Bias² by subtracting the two variance estimates from the estimated loss.

For each of these three quantities, Bias², PretVar and FineVar, we estimate it for each instance, average it across all instances in the evaluation set, and report it in Table 7. The variances at the instance level are much larger than the variance of overall accuracy, by a factor of 1000.

We may conclude from Table 7 that larger models have larger finetuning variance and smaller pretraining variance. However, lower bias also inherently implies lower variance. To see this, suppose a model has perfect accuracy and hence zero bias;

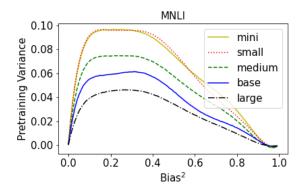


Figure 5: The pretraining variance conditioned on Bias² (the level of correctness). Each curve represents a model size. Larger models have lower pretraining variance across all levels of bias.

then it always predicts the same label (the correct one) and hence has zero variance. This might favor larger models and "underestimate" their variance, since they have lower bias. Therefore, we calculate and compare the variances conditioned on the bias, i.e. $\operatorname{PretVar}(b^2) := \mathbb{E}_i[\operatorname{PretVar}_i|\operatorname{Bias}^2_i = b^2]$.

We estimate $\operatorname{PretVar}^s(b^2)$ using Gaussian process regression and plot it against b^2 in Figure 5. We find that larger models still have lower pretraining variance across all levels of bias on the specific task of MNLI under the 0/1 loss. To further check whether our conclusions are general, we tested them on other tasks and under the squared loss $\mathcal{L}_i := (1-p_i)^2$, where p_i is the probability assigned to the correct class. Below are the conclusions that generally hold across different tasks and loss functions.

Conclusion We find that 1) larger models have larger finetuning variance, 2) LARGE has smaller pretraining variance than BASE; however, the ordering between other sizes varies across tasks and losses, and 3) finetuning variance is 2–8 times as large as pretraining variance, and the ratio is bigger for larger models.

7 Discussion and Future Directions

To investigate model behaviors at the instance level, we produced massive amounts of model predictions in Section 2 and treated them as raw data. To extract insights from them, we developed better metrics and statistical tools, including a new method to control the false discoveries, an unbiased estimator for the decomposed variances, and metrics that compute variance and correlation of improvements

conditioned on instance accuracy. We find that larger pretrained models are indeed worse on a non-trivial fraction of instances and have higher variance due to finetuning seeds; additionally, instance accuracy improvements from MINI to MEDIUM correlate with improvements from MEDIUM to LARGE

Overall, we treated model prediction data as the central object and built analysis tools around them to obtain a finer understanding of model performance. We therefore refer to this paradigm as "instance level understanding as data mining". We discuss three key factors for this paradigm to thrive: 1) scalability and the cost of obtaining prediction data, 2) other information to collect for each instance, and 3) better statistical tools. We analyze each of these aspects below.

Scalability and Cost of Data Data mining is more powerful with more data. How easy is it to obtain more model predictions? In our paper, the main bottleneck is pretraining. However, once the pretrained models are released, individual researchers can download them and only need to repeat the cheaper finetuning procedure.

Furthermore, model prediction data are undershared: while many recent research papers share their code or even model weights to help reproduce the results, it is not yet a standard practice to share all the model predictions. Since many researches follow almost the same recipe of pretraining and finetuning (McCoy et al., 2020; Desai and Durrett, 2020; Dodge et al., 2020), much computation can be saved if model predictions are shared. On the other hand, as the state of the art model size is increasing at a staggering speed⁷, most researchers will not be able to run inference on a single instance. The trend that models are becoming larger and more similar necessitate more prediction sharing.

Meta-Labels and Other Predictions Data mining is more powerful with more types of information. One way to add information to each instance is to assign "meta-labels". In the HANS (McCoy et al., 2019) dataset, the authors tag each instance with a heuristic ⁸ that holds for the training distribution but fails on this instance. Naik et al. (2018a) and Ribeiro et al. (2020) associate each instance

⁷e.g. BERT (Devlin et al., 2019) has 340M parameters, while Switch-Transformer has over 1 trillion parameters (Fedus et al., 2021).

⁸For example, "the label [entailment] is likely if the premise and the hypothesis have significant lexical overlap".

with a particular stress test type or subgroup, for example, whether the instance requires the model to reason numerically or handle negations. Nie et al. (2020) collects multiple human responses to estimate human disagreement for each instance. This meta-information can potentially help us identify interpretable patterns for the disagreeing instances where one model is better than the other. On the flip side, identifying disagreeing instances between two models can also help us generate hypothesis and decide what subgroup information to annotate.

We can also add performance information on other tasks to each instance. For example, Pruksachatkun et al. (2020) studied the correlation between syntactic probing accuracy (Hewitt and Liang, 2019) and downstream task performance. Turc et al. (2019) and Kaplan et al. (2020) studied the correlation between language modelling loss and the downstream task performance. However, they did not analyze correlations at the instance level. We may investigate whether their results hold on the instance level: if an instance is easier to tag by a probe or easier to predict by a larger language model, is the accuracy likely to be higher?

Statistical Tools Data mining is more powerful with better statistical tools. Initially we used the Benjamini-Hochberg procedure with Fisher's exact test, which required us to pretrain 10 models to formally verify that the decaying instances exist. However, we later realized that 2 is in fact enough by using our approach introduced in Section 4. We could have saved 80% of the computation for pretraining if this approach was known before we started.

Future work can explore more complicated metrics and settings. We compared at most 3 different model sizes at a time, and higher order comparisons require novel metrics. We studied two sources of randomness, pretraining and finetuning, but other sources of variation can be interesting as well, e.g. differences in pretraining corpus, different model checkpoints, etc. To deal with more sophisticated metrics, handle different sources and hierarchies of randomness, and reach conclusions that are robust to noises at the instance level, researchers need to develop new inference procedures.

To conclude, for better instance level understanding, we need to produce and share more prediction data, annotate more diverse linguistic properties, and develop better statistical tools to infer under noises. We hope our work can inform researchers

about the core challenges underlying instance level understanding and inspire future work.

Acknowledgement

We thank Steven Cao, Cathy Chen, Frances Ding, David Gaddy, Colin Li, and Alex Wei for giving comments on the initial paper draft. We would also like to thank the Google Cloud TPU team for their hardware support.

References

Yoav Benjamini and Yosef Hochberg. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300.

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

Alexander D'Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D Hoffman, et al. 2020. Underspecification presents challenges for credibility in modern machine learning. arXiv preprint arXiv:2011.03395.

Shrey Desai and Greg Durrett. 2020. Calibration of pre-trained transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 295–302, Online. Association for Computational Linguistics.

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.

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. arXiv preprint arXiv:2002.06305.

William Fedus, Barret Zoph, and Noam Shazeer. 2021. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *arXiv* preprint arXiv:2101.03961.

- Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. 2020. Pretrained transformers improve out-of-distribution robustness. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2744–2751, Online. Association for Computational Linguistics.
- Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. 2020. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*.
- John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. arXiv preprint arXiv:2001.08361.
- Wuwei Lan, Siyu Qiu, Hua He, and Wei Xu. 2017. A continuously growing dataset of sentential paraphrases. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1224–1234, Copenhagen, Denmark. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer.
 2020. BART: Denoising sequence-to-sequence pretraining for natural language generation, translation, and comprehension. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7871–7880, Online. Association for Computational Linguistics.
- Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, Kurt Keutzer, Dan Klein, and Joseph E. Gonzalez. 2020. Train large, then compress: Rethinking model size for efficient training and inference of transformers. In *International Conference on Machine Learn*ing.
- 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 *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 217–227, Online. Association for Computational Linguistics.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association*

- for Computational Linguistics, pages 3428–3448, Florence, Italy. Association for Computational Linguistics.
- Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018a. Stress test evaluation for natural language inference. In *The 27th International Conference on Computational Linguistics (COLING)*, Santa Fe, New Mexico, USA.
- Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018b. Stress test evaluation for natural language inference. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Yixin Nie, Xiang Zhou, and Mohit Bansal. 2020. What can we learn from collective human opinions on natural language inference data? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9131–9143, Online. Association for Computational Linguistics.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Yada Pruksachatkun, Jason Phang, Haokun Liu, Phu Mon Htut, Xiaoyi Zhang, Richard Yuanzhe Pang, Clara Vania, Katharina Kann, and Samuel R. Bowman. 2020. Intermediate-task transfer learning with pretrained language models: When and why does it work? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5231–5247, Online. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 4902–4912, Online. Association for Computational Linguistics.

- Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. 2020. An investigation of why overparameterization exacerbates spurious correlations. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 8346–8356. PMLR.
- Lane Schwartz, Timothy Anderson, Jeremy Gwinnup, and Katherine Young. 2014. Machine translation and monolingual postediting: The AFRL WMT-14 system. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 186–194, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962v2*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Adina Williams, Tiago Pimentel, Hagen Blix, Arya D. McCarthy, Eleanor Chodroff, and Ryan Cotterell. 2020. Predicting declension class from form and meaning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6682–6695, Online. Association for Computational Linguistics.

A Pretraining and Finetuning Details

Here we explain how to obtain the model predictions, which are analyzed in later sections. To obtain these predictions under the "pretraining and finetuning" framework (Devlin et al., 2019), we need to decide a model size, perform pretraining, finetune on a training set with a choice of hyperparameters, and test the model on an evaluation set. We discuss each bolded aspects below.

Size Similar to Turc et al. (2019), we experimented with the following five model sizes, listed in increasing order: MINI (L4/H256) 9 , SMALL (L4/H512), MEDIUM (L8/H512), BASE (L12/H768), and LARGE (L24/H1024).

Pretraining We used the pretraining code from Devlin et al. (2019) and the pre-training corpus from Li et al. (2020). Compared to the original BERT release, we used context size 128 instead of 512, since computation cost grows quadratically with respect to context size; we also pretrained for 2M steps instead of 1M.

Training Set We consider 3 datasets: Quora Question Pairs (QQP) ¹⁰, Multi-Genre Natural Language Inference (MNLI; Williams et al. (2020)), and the Stanford Sentiment Treebank (SST-2; (Socher et al., 2013)). For QQP we used the official training split. For MNLI we used 350K out of 400K instances from the original training split, and added the remaining 50K to the evaluation set, since the original in-domain development set only contains 10K examples. For SST-2, we mix the training and development set of the original split, split the instances into 5 folds, train on four of them, and evaluate on the remaining fold.

Hyperparameters As in Turc et al. (2019), we finetune 4 epochs for each dataset. For each task and model size, we tune hyperparameters in the following way: we first randomly split our new training set into 80% and 20%; then we finetune on the 80% split with all 9 combination of batch size [16, 32, 64] and learning rate [1e-4, 5e-5, 3e-5], and choose the combination that leads to the best average accuracy on the remaining 20%.

Evaluation Set After finetuning our pretrained models, we evaluate them on a range of in-domain,

out-of-domain, or challenging datasets to obtain model predictions. Models trained on MNLI are also evaluated on Stanford Natural Language Inference (SNLI; (Bowman et al., 2015)), Heuristic Analysis for NLI Systems (HANS; (McCoy et al., 2019)), and stress test evaluations (STRESS; (Naik et al., 2018b)). Models trained on QQP are also evaluated on Twitter Paraphrase Database (TwitterPPDB; (Lan et al., 2017)).

Since pretraining introduces randomness, for each model size s, we pretrain 10 times with different random seed P; since finetuning also introduces noise, for each pretrained model we pretrain 5 times with different random seed F; besides, we also evaluate the model at the checkpoints after E epochs, where $E \in [3, 3\frac{1}{3}, 3\frac{2}{3}, 4]$.

Pretraining 10 models for all 5 model sizes altogether takes around 3840 hours on TPU v3 with 8 cores. Finetuning all of them 5 times for all three tasks in our paper requires around 1200 hours.

B Compare Our Models to the Original

Since we decreased the pre-training context length to save computation, these models are not exactly the same as the original BERT release by Devlin et al. (2019) and Turc et al. (2019). We need to benchmark our model against theirs to ensure that the performance of our model is still reasonable and the qualitative trend still holds. For each each size and task, we finetune the original model 5 times and calculate the average of overall accuracy.

The comparison can be seen in Table 8. We find that our model does not substantially differ from the original ones on QQP and SST-2. On MNLI, the performance of our BERT-BASE and BERT-LARGE is $2\sim3\%$ below the original release, but the qualitative trend that larger models have better accuracy still holds robustly.

C More Instance Difference Results

Similar to Figure 4, for all 10 pairs of model sizes and all in-distribution instances of MNLI, SST-2, and QQP, we plot the cumulative density of $\Delta \hat{A}cc$ and $\Delta Acc'$, or say, $\hat{D}ecay(t)$ and $\hat{D}ecay'(t)$ in Figure 6, 7, and 8.

Additionally, for each pair of model sizes s_1 and s_2 , we estimate "how much instances are getting better/worse accuracy?" by taking the maximum difference between the red curve and the blue curve. We report these results for MNLI, SST-2, and QQP in Table 9. We find that larger model size gaps

⁹4 Layers with hidden dimension 256

¹⁰https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs

	QQP	MNLI	SST-2
MINI ^{orig}	88.2%	74.6%	92.8%
MINI ours	87.3%	74.3%	92.8%
SMALL orig	89.1%	77.3%	93.9%
SMALL ours	88.7%	76.7%	93.9%
MEDIUM ^{orig}	89.8%	79.6%	94.2%
MEDIUM ours	89.5%	78.9%	94.2%
BASE ^{orig}	90.8%	83.8%	95.0%
BASE ours	90.6%	81.2%	94.6%
LARGE ^{orig}	91.3%	86.8%	95.2%
LARGE ours	91.0%	83.8%	94.8%

Table 8: Comparing our pretrained model (superscript orig) to the original release by Devlin et al. (2019) and Turc et al. (2019) (superscript ours). All pretrained models are finetuned with the training set and tested on the in-distribution evaluation set described in Appendix A.

lead to larger decaying fraction, but also larger improving fraction as well.

D Proof of Theorem 1

Formal Setup Our goal is to show that if all the random seeds are independent,

$$Decay > \mathbb{E}[\hat{Decay}(t) - Decay'(t)]$$
 (16)

More concretely, suppose each instance is indexed by i, the set of all instances is \mathcal{T} , and the random seed is R; then $c_R^s \in \{0,1\}^{|\mathcal{T}|}$ is a random $|\mathcal{T}|$ dimensional vector, where $c_{R,i}^s = 1$ if the model of size s correctly predicts instance i under the random seed R. We are comparing model size s_1 and s_2 , where s_2 is larger; to keep notation uncluttered, we omit these indexes whenever possible.

Suppose we observe $c_{R_1...2k}^{s_1}$ and $c_{R_{2k+1...4k}}^{s_2}$, where there are 2k different random seeds for each model size 11 . Then

$$\Delta \hat{A} cc_i := \frac{1}{2k} \left(\sum_{j=1}^{2k} c_{R_j,i}^{s_1} - \sum_{j=2k+1}^{4k} c_{R_j,i}^{s_2} \right), \quad (17)$$

and hence the discovery rate Decay(t) is defined as

$$\hat{\text{Decay}}(t) := \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} \mathbf{1}[\Delta \hat{\text{Acc}} \le t].$$
 (18)

For the random baseline estimator, we have

$$\Delta \operatorname{Acc}_{i}' := \frac{1}{2k} \left(\sum_{j=1}^{k} c_{R_{j},i}^{s_{1}} + \sum_{j=2k+1}^{3k} c_{R_{j},i}^{s_{2}} - \sum_{j=k+1}^{2k} c_{R_{j},i}^{s_{1}} - \sum_{j=3k+1}^{4k} c_{R_{j},i}^{s_{2}} \right),$$
(19)

and the false discovery control Decay' is defined as

$$Decay'(t) := \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} \mathbf{1}[\Delta Acc_i' \le t].$$
 (20)

To reiterate, the definition of the true decay rate is

$$Decay = \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} \mathbf{1}[\Delta Acc_i < 0].$$
 (21)

Our goal is to prove that

$$\operatorname{Decay} \geq \mathbb{E}_{R_1...R_{4k}}[\hat{\operatorname{Decay}}(t) - \operatorname{Decay}'(t)]$$
 (22)

Proof By re-arranging terms and linearity of expectation, Equation 22 is equivalent to the following

$$\sum_{i=1}^{|\mathcal{T}|} (\mathbf{1}[\Delta A c c_i < 0] - \mathbb{P}[\Delta \hat{A} c c_i \le t] + \mathbb{P}[\Delta A c c_i' \le t]) \ge 0$$
(23)

Hence, we can declare victory if we can prove that for all i,

$$\mathbf{1}[\Delta Acc_i < 0] - \mathbb{P}[\Delta \hat{A}cc_i \le t]$$
 (24)

$$+ \mathbb{P}[\Delta Acc_i' \le t] \ge 0$$

To prove Equation 24, we observe that if $\mathrm{Acc}_i < 0$, since the probabilities are bounded by 0 and 1, its left-hand side must be positive. Therefore, we only need to prove that

$$\Delta \operatorname{Acc}_{i} \ge 0$$

$$\Rightarrow \mathbb{P}[\Delta \operatorname{Acc}'_{i} \le t] \ge \mathbb{P}[\Delta \operatorname{Acc}_{i} \le t],$$
(25)

which will be proved in Lemma 1.□

Lemma 1

$$\Delta Acc_i \ge 0$$

$$\Rightarrow \mathbb{P}[\Delta Acc_i' \le t] \ge \mathbb{P}[\Delta \hat{A}cc_i \le t],$$
(26)

¹¹We assumed even number of random seeds since we will mix half of the models from each size to compute the random baseline

MNLI $s_1 \setminus s_2$	MINI	SMALL	MEDIUM	BASE	LARGE
MINI	0.000	0.087	0.136	0.179	0.214
SMALL	0.033	0.000	0.089	0.139	0.180
MEDIUM	0.050	0.028	0.000	0.090	0.143
BASE	0.060	0.048	0.026	0.000	0.101
LARGE	0.059	0.052	0.040	0.021	0.000
$QQP s_1 \setminus s_2$	MINI	SMALL	MEDIUM	BASE	LARGE
MINI	0.000	0.057	0.076	0.100	0.107
SMALL	0.019	0.000	0.039	0.073	0.084
MEDIUM	0.029	0.014	0.000	0.044	0.063
BASE	0.034	0.027	0.016	0.000	0.032
LARGE	0.036	0.031	0.027	0.016	0.000
SST-2 $s_1 \setminus s_2$	MINI	SMALL	MEDIUM	BASE	LARGE
MINI	0.000	0.037	0.043	0.052	0.057
SMALL	0.010	0.000	0.015	0.031	0.036
MEDIUM	0.016	0.008	0.000	0.020	0.028
BASE	0.019	0.014	0.009	0.000	0.014
LARGE	0.020	0.017	0.015	0.008	0.000

Table 9: On QQP, MNLI in domain development set and SST-2 we lowerbound the fraction of instances that improves when model size changes from s_1 (row) to s_2 (column).

For m = 1, 2, define

$$p_i^{s_m} := \mathbb{E}_R[c_i^{s_m}],\tag{27}$$

then

$$p_i^{s_1} \le p_i^{s_2} \tag{28}$$

Since $c_i^{s_1}$ and $c_i^{s_2}$ are both Bernoulli random variables with rate $p_i^{s_1}$ and $p_i^{s_2}$ respectively, we can write down the probability distribution of $\Delta \hat{A}cc_i$ and $\Delta Acc_i'$ as the sum/difference of several binomial variables, i.e.

$$\begin{split} &\Delta \hat{\mathbf{A}} \mathbf{cc}_i \sim (\mathsf{Binom}(k, p_i^{s_2}) + \mathsf{Binom}(k, p_i^{s_2}) \ \ (29) \\ &- \mathsf{Binom}(k, p_i^{s_1}) - \mathsf{Binom}(k, p_i^{s_1}))/2k, \end{split}$$

and

$$\Delta \text{Acc}'^i \sim (\text{Binom}(k, p^{s_2, i}) + \text{Binom}(k, p^{s_1, i})$$
 (30)

$$-\operatorname{Binom}(k,p^{s_1,i})-\operatorname{Binom}(k,p^{s_2,i}))/2k$$

 $p_i^{s_1} \leq p_i^{s_2}$, $\operatorname{Binom}(k, p^{s_2,i})$) first order stochastically dominates $\operatorname{Binom}(k, p^{s_1,i})$. Therefore, $\Delta \operatorname{Acc}^{\prime i}$ dominates $\Delta \operatorname{Acc}_i$, hence completing the proof. \Box

D.1 Independent Seed Assumption

We notice that Theorem 1 requires the seeds R to be independent. This assumption does not hold on

our data, since some finetuning runs share the same pretraining seeds. Therefore, the above proof no longer holds. Specifically, Lemma 1 fails because $\Delta \hat{A} cc$ and $\Delta Acc'$ are no longer binomial variables, and the later does not necessarily dominate the first. Here is a counter-example, if the seeds are not entirely independent.

Hypothetically, suppose we are comparing a smaller model s_1 and a larger model s_2 . For the smaller model, with .1 probability it finds a perfect pretrained model that always predict correctly across all finetuning runs and with .9 probability it finds a bad pretrained model that predict always incorrectly. For the larger model, with probability 1 it finds an average pretrained model that predict correctly for .2 fraction of finetuning runs. The larger model is on average better, because it has .2 > .1 probability to be correct. Hence, $\Delta {\rm Acc} > 0$

Suppose we observe 2 independent pretraining seeds for each size and infinite number of fine-tuning seeds for each pretraining seed, and let us consider the threshold -0.8. Then

$$\mathbb{P}[\Delta \hat{A} cc_i \le -0.8] \qquad (31)$$

$$= 0.01 \ge 0 = \mathbb{P}[\Delta Acc_i' \le -0.8]$$
 (32)

The event that $\Delta \hat{A}cc_i \leq -0.8$ happens with probability 0.01 when both of the two small pretrained models have good pretraining seeds, and $\Delta Acc_i'$ is at least -0.5 and will never be less than -0.8.

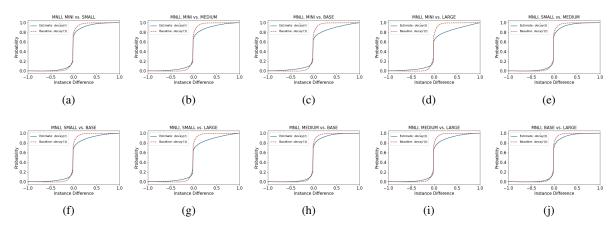


Figure 6: Similar to Figure 4, on MNLI in-distribution development set, for each pair of model sizes, we plot the cumulative density function of instance differences.

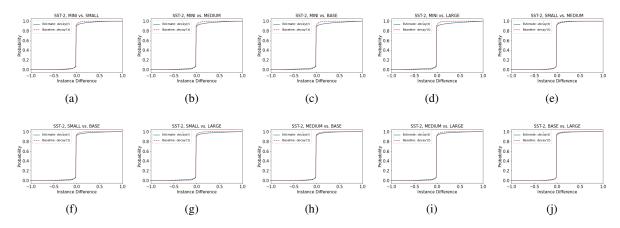


Figure 7: Similar to Figure 4, on SST-2, for each pair of model sizes, we plot the cumulative density function of instance differences.

The key idea behind this counter-example is that even if the larger model has better average, the distribution of average finetuning accuracy for different pretraining seeds might not stochastically dominate the one with lower average because of outliers. Hence, a priori, this is unlikely to happen in practice, since pretraining variance is generally small, and we have multiple pretraining seeds to average out the outliers. Nevertheless, future work is needed to make a more rigorous argument.

E Upward Bias of Adaptive Thresholds

In section 3 we picked the best threshold that can maximize the lowerbound, which can incur a slight upward bias. Here we estimate that the bias is at most 10% relative to the unbiased lowerbound with a bootstrapping method.

We use the empirical distribution of 10 pretrained models as the ground truth distribution for bootstrapping. We first compute a best threshold with 10 sampled smaller and larger pretrained models, and then compute the lowerbound L with this threshold on another sample of 10 smaller and larger models. Intuitively, we use one bootstrap sample (which contains 10 smaller pretrained models and 10 larger pretrained models) as the development set to "tune the threshold", and then use this threshold on a fresh bootstrap sample to compute the lowerbound. We refer to the lowerbound that uses the best threshold as L^* , and compute the relative error $\mathbb{E}[(L^*-L)]/\mathbb{E}[L)]$, where the expectation is taken with respect to bootstrap samples.

We report all results in Table 10. In general, we find that the upward bias is negligible, which is at most around 10%.

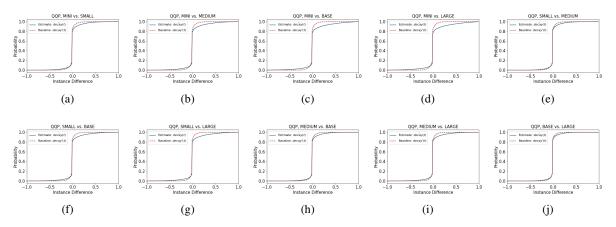


Figure 8: Similar to Figure 4, on QQP in-domain development set, for each pair of model sizes, we plot the cumulative density function of instance differences.

F Comparison with Significance Testing

We also experimented with the classical approach that calculates the significance-level for each instance and then use the Benjamini-Hochberg procedure to lowerbound the decaying fraction. To make sure that we are comparing with this approach fairly, we lend it additional power by picking the false discovery rate that can maximize the true discovery counts. We report the decaying fraction on MNLI in-domain development set found by this classical method and compare it with our method for different model size differences in Table 11; we also simulate situations when we have fewer models.

In general, we find that our method always provide a tighter (higher) lowerbound than the classical method, and 2 models are sufficient to verify the existence (i.e. lowerbound > 0) of the decaying fraction; in contrast, the classical method sometimes fails to do this even with 10 models, e.g., when comparing BASE to LARGE.

Intuitively, our approach provides a better lower-bound because it better makes use of the information that on most instances, both the smaller and the larger models agree and predict completely correctly or incorrectly ¹². For an extreme example, suppose we only observe 2 smaller models and 2 larger models, and infinite number of datapoints, whose predictions are independent. On 99.98% datapoints, both models have instance accuracy 1; on 0.01% datapoints, smaller model is completely correct while bigger completely wrong,

while on the rest 0.01% smaller completely wrong but bigger completely correct. Setting threshold to be 2, our decay estimate Decay is 0.01%, while Decay' = 0: since the models either completely predict correct or wrongly, there is never a false discovery. Therefore, our method can provide the tightest lowerbound 0.01% in this case. On the other hand, since we only have 4 models in total, the lowest significance-level given by the fisher exact test is $17\% \gg 0.1\%$, hence the discovery made by the Benjamin-Hochberg procedure is 0.

G More Results on Momentum

We report more results on the correlation between instance differences. Specifically, for one triplet of model sizes (e.g. MINI \Rightarrow MEDIUM \Rightarrow LARGE), for each group of instances that have similar $\hat{\mathrm{Acc}}^{\mathrm{MEDIUM}}$, we calculate the correlation between instance differences, i.e. the Pearson-R score between $_{\mathrm{MEDIUM}}^{\mathrm{MINI}}\Delta\mathrm{Acc}$ and $_{\mathrm{LARGE}}^{\mathrm{MEDIUM}}\Delta\mathrm{Acc}$. All results can be seen in Table 12.

We observe that

- For nearly all buckets, the improvements are positively correlated.
- When model size gap becomes larger (e.g. MINI, MEDIUM, LARGE has the largest model size differences), the correlation decreases.

H Loss Decomposition and Estimation

In this section, under the bias-variance decomposition and total variance decomposition framework, we decompose loss into four components: bias, variance brought by pretraining randomness, by

¹²This is for intuition, though, and we do not need any assumption on the prior of instance accuracy, which requires a Bayes interpretation.

MNLI $s_1 \setminus s_2$	MINI	SMALL	MEDIUM	BASE	LARGE
MINI	0.000	0.031	0.027	0.026	0.020
SMALL	0.108	0.000	0.027	0.023	0.019
MEDIUM	0.095	0.116	0.000	0.028	0.023
BASE	0.093	0.100	0.144	0.000	0.026
LARGE	0.097	0.103	0.117	0.149	0.000
$QQP s_1 \setminus s_2$	MINI	SMALL	MEDIUM	BASE	LARGE
MINI	0.000	0.025	0.022	0.021	0.020
SMALL	0.127	0.000	0.040	0.020	0.020
MEDIUM	0.093	0.146	0.000	0.032	0.031
BASE	0.087	0.119	0.123	0.000	0.049
LARGE	0.090	0.105	0.079	0.106	0.000
SST-2 $s_1 \setminus s_2$	MINI	SMALL	MEDIUM	BASE	LARGE
MINI	0.000	0.028	0.022	0.021	0.019
SMALL	0.117	0.000	0.047	0.031	0.029
MEDIUM	0.075	0.093	0.000	0.063	0.035
BASE	0.068	0.067	0.085	0.000	0.071
LARGE	0.071	0.067	0.060	0.098	0.000

Table 10: The same table as 10, except that we are now calculating the relative upward bias $\mathbb{E}[(L^*-L)]/\mathbb{E}[L)]$ as described in Section E.

s1	s2	method	2	4	6	8	10
MINI	SMALL	ours	0.004	0.011	0.016	0.020	0.023
MINI	SMALL	BH	0.000	0.000	0.000	0.000	0.002
MINI	MEDIUM	ours	0.012	0.019	0.026	0.032	0.035
MINI	MEDIUM	BH	0.000	0.000	0.003	0.008	0.011
MINI	BASE	ours	0.019	0.028	0.035	0.040	0.042
MINI	BASE	BH	0.000	0.000	0.008	0.014	0.020
MINI	LARGE	ours	0.019	0.027	0.031	0.037	0.040
MINI	LARGE	BH	0.000	0.000	0.009	0.015	0.019
SMALL	MEDIUM	ours	0.002	0.006	0.010	0.015	0.017
SMALL	MEDIUM	BH	0.000	0.000	0.000	0.000	0.000
SMALL	BASE	ours	0.013	0.020	0.025	0.030	0.033
SMALL	BASE	BH	0.000	0.000	0.002	0.006	0.011
SMALL	LARGE	ours	0.015	0.021	0.026	0.031	0.033
SMALL	LARGE	BH	0.000	0.000	0.005	0.009	0.013
MEDIUM	BASE	ours	0.006	0.010	0.013	0.014	0.016
MEDIUM	BASE	BH	0.000	0.000	0.000	0.000	0.001
MEDIUM	LARGE	ours	0.010	0.014	0.019	0.022	0.023
MEDIUM	LARGE	BH	0.000	0.000	0.002	0.004	0.006
BASE	LARGE	ours	0.004	0.005	0.009	0.010	0.012
BASE	LARGE	BH	0.000	0.000	0.000	0.000	0.000

Table 11: We compare each pair of model sizes s_1 and s_2 and report the lower bound provided by our method and the Benjamin-Hochberg (BH) procedure. The numbers in column name denote how many pretrained model we used to obtain the lower bounds.

0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	1.00
0.00	0.18	0.19	0.18	0.23	0.26	0.24	0.23	0.20	0.12
0.07	0.22	0.29	0.40	0.35	0.33	0.38	0.27	0.24	0.13
0.05	0.09	0.17	0.33	0.20	0.30	0.12	0.13	0.16	0.09
0.03	0.15	0.18	0.33	0.17	0.16	0.22	0.20	0.19	0.09
0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	1.00
0.03	0.21	0.18	0.21	0.21	0.25	0.18	0.16	0.10	0.06
0.01	0.17	0.23	0.19	0.24	0.22	0.24	0.19	0.16	0.05
-0.02	0.16	0.09	0.23	0.17	0.10	0.14	0.14	0.09	-0.01
-0.01	0.07	0.14	0.09	0.16	0.09	0.16	0.07	0.10	0.07
0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	1.00
0.09	0.26	0.43	0.22	0.28	0.24	0.27	0.35	0.20	0.12
0.07	0.12	0.22	0.40	0.07	0.20	0.10	0.12	0.19	0.06
0.01	0.24	0.29	0.35	0.19	0.19	0.26	0.39	0.15	0.03
0.01	0.17	0.11	0.41	0.04	0.29	0.16	0.21	0.15	0.07
	0.00 0.07 0.05 0.03 0.10 0.03 0.01 -0.02 -0.01 0.10 0.09 0.07 0.01	0.00 0.18 0.07 0.22 0.05 0.09 0.03 0.15 0.10 0.20 0.03 0.21 0.01 0.17 -0.02 0.16 -0.01 0.07 0.10 0.20 0.09 0.26 0.07 0.12 0.01 0.24	0.00 0.18 0.19 0.07 0.22 0.29 0.05 0.09 0.17 0.03 0.15 0.18 0.10 0.20 0.30 0.03 0.21 0.18 0.01 0.17 0.23 -0.02 0.16 0.09 -0.01 0.07 0.14 0.10 0.20 0.30 0.09 0.26 0.43 0.07 0.12 0.22 0.01 0.24 0.29	0.00 0.18 0.19 0.18 0.07 0.22 0.29 0.40 0.05 0.09 0.17 0.33 0.03 0.15 0.18 0.33 0.10 0.20 0.30 0.40 0.03 0.21 0.18 0.21 0.01 0.17 0.23 0.19 -0.02 0.16 0.09 0.23 -0.01 0.07 0.14 0.09 0.10 0.20 0.30 0.40 0.09 0.26 0.43 0.22 0.07 0.12 0.22 0.40 0.01 0.24 0.29 0.35	0.00 0.18 0.19 0.18 0.23 0.07 0.22 0.29 0.40 0.35 0.05 0.09 0.17 0.33 0.20 0.03 0.15 0.18 0.33 0.17 0.10 0.20 0.30 0.40 0.50 0.03 0.21 0.18 0.21 0.21 0.01 0.17 0.23 0.19 0.24 -0.02 0.16 0.09 0.23 0.17 -0.01 0.07 0.14 0.09 0.16 0.10 0.20 0.30 0.40 0.50 0.09 0.26 0.43 0.22 0.28 0.07 0.12 0.22 0.40 0.07 0.01 0.24 0.29 0.35 0.19	0.00 0.18 0.19 0.18 0.23 0.26 0.07 0.22 0.29 0.40 0.35 0.33 0.05 0.09 0.17 0.33 0.20 0.30 0.03 0.15 0.18 0.33 0.17 0.16 0.10 0.20 0.30 0.40 0.50 0.60 0.03 0.21 0.18 0.21 0.21 0.25 0.01 0.17 0.23 0.19 0.24 0.22 -0.02 0.16 0.09 0.23 0.17 0.10 -0.01 0.07 0.14 0.09 0.16 0.09 0.10 0.20 0.30 0.40 0.50 0.60 0.09 0.26 0.43 0.22 0.28 0.24 0.07 0.12 0.22 0.40 0.07 0.20 0.01 0.24 0.29 0.35 0.19 0.19	0.00 0.18 0.19 0.18 0.23 0.26 0.24 0.07 0.22 0.29 0.40 0.35 0.33 0.38 0.05 0.09 0.17 0.33 0.20 0.30 0.12 0.03 0.15 0.18 0.33 0.17 0.16 0.22 0.10 0.20 0.30 0.40 0.50 0.60 0.70 0.03 0.21 0.18 0.21 0.21 0.25 0.18 0.01 0.17 0.23 0.19 0.24 0.22 0.24 -0.02 0.16 0.09 0.23 0.17 0.10 0.14 -0.01 0.07 0.14 0.09 0.16 0.09 0.16 0.10 0.20 0.30 0.40 0.50 0.60 0.70 0.09 0.26 0.43 0.22 0.28 0.24 0.27 0.07 0.12 0.22 0.40 0.07 0.20 <td>0.00 0.18 0.19 0.18 0.23 0.26 0.24 0.23 0.07 0.22 0.29 0.40 0.35 0.33 0.38 0.27 0.05 0.09 0.17 0.33 0.20 0.30 0.12 0.13 0.03 0.15 0.18 0.33 0.17 0.16 0.22 0.20 0.10 0.20 0.30 0.40 0.50 0.60 0.70 0.80 0.03 0.21 0.18 0.21 0.21 0.25 0.18 0.16 0.01 0.17 0.23 0.19 0.24 0.22 0.24 0.19 -0.02 0.16 0.09 0.23 0.17 0.10 0.14 0.14 -0.01 0.07 0.14 0.09 0.16 0.09 0.16 0.07 0.10 0.20 0.30 0.40 0.50 0.60 0.70 0.80 0.09 0.26 0.43 0.22<</td> <td>0.00 0.18 0.19 0.18 0.23 0.26 0.24 0.23 0.20 0.07 0.22 0.29 0.40 0.35 0.33 0.38 0.27 0.24 0.05 0.09 0.17 0.33 0.20 0.30 0.12 0.13 0.16 0.03 0.15 0.18 0.33 0.17 0.16 0.22 0.20 0.19 0.10 0.20 0.30 0.40 0.50 0.60 0.70 0.80 0.90 0.03 0.21 0.18 0.21 0.21 0.25 0.18 0.16 0.10 0.01 0.17 0.23 0.19 0.24 0.22 0.24 0.19 0.16 -0.02 0.16 0.09 0.23 0.17 0.10 0.14 0.09 -0.01 0.07 0.14 0.09 0.16 0.09 0.16 0.07 0.10 0.10 0.20 0.30 0.40 0.</td>	0.00 0.18 0.19 0.18 0.23 0.26 0.24 0.23 0.07 0.22 0.29 0.40 0.35 0.33 0.38 0.27 0.05 0.09 0.17 0.33 0.20 0.30 0.12 0.13 0.03 0.15 0.18 0.33 0.17 0.16 0.22 0.20 0.10 0.20 0.30 0.40 0.50 0.60 0.70 0.80 0.03 0.21 0.18 0.21 0.21 0.25 0.18 0.16 0.01 0.17 0.23 0.19 0.24 0.22 0.24 0.19 -0.02 0.16 0.09 0.23 0.17 0.10 0.14 0.14 -0.01 0.07 0.14 0.09 0.16 0.09 0.16 0.07 0.10 0.20 0.30 0.40 0.50 0.60 0.70 0.80 0.09 0.26 0.43 0.22<	0.00 0.18 0.19 0.18 0.23 0.26 0.24 0.23 0.20 0.07 0.22 0.29 0.40 0.35 0.33 0.38 0.27 0.24 0.05 0.09 0.17 0.33 0.20 0.30 0.12 0.13 0.16 0.03 0.15 0.18 0.33 0.17 0.16 0.22 0.20 0.19 0.10 0.20 0.30 0.40 0.50 0.60 0.70 0.80 0.90 0.03 0.21 0.18 0.21 0.21 0.25 0.18 0.16 0.10 0.01 0.17 0.23 0.19 0.24 0.22 0.24 0.19 0.16 -0.02 0.16 0.09 0.23 0.17 0.10 0.14 0.09 -0.01 0.07 0.14 0.09 0.16 0.09 0.16 0.07 0.10 0.10 0.20 0.30 0.40 0.

Table 12: Sorted in ascending order, the model sizes are MINI , SMALL , MEDIUM , BASE , and LARGE . The three model sizes listed for each row represents the model size of interest: for example, MINI, MEDIUM, LARGE means that we are calculating the correlation between $_{\text{MEDIUM}}^{\text{MINI}}\Delta\text{Acc}$ and $_{\text{LARGE}}^{\text{MEDIUM}}\Delta\text{Acc}$. Each column t represents a bucket that contains instances with middle size accuracy in [t-0.1,t]. For example, if the row name is MINI, MEDIUM, LARGE, then the column 0.2 corresponds to a bucket where $\hat{\text{Acc}}_{\hat{k}}^{\text{MEDIUM}}$ is between 0.1 and 0.2. We calculate the PearsonR correlation score between $_{\text{MEDIUM}}^{\text{MINI}}\Delta\hat{\text{Acc}}$ and $_{\text{LARGE}}^{\text{MEDIUM}}\Delta\hat{\text{Acc}}$ across all instances in the bucket.

finetuning randomness, and across different checkpoints throughout training. We formally define the quantities we want to estimate in Appendix H.1, present an unbiased estimator for these quantities in Appendix H.2, and show that our method can be generalized to arbitrary number of source of randomness in Appendix H.3.

Specifically, the main paper focused on scenarios with 2 sources of randomness: pretraining and finetuning. We discuss the case with 3 sources of randomness in the appendix, rather than 2 as in the main paper, because it is easier to understand the general estimation strategy in the case of 3.

H.1 Formalizing Decomposition

Recall that P is the pretraining seed, F is the fine-tuning seed, E represents a model checkpoint, i indexes each instance (datapoint). $c_{P,F,E}^{s,i}=1$ if the model of size s with pretraining seed p and fine-tuning seed F, and trained for E epochs is correct on datapoint i, and 0 otherwise. Notice that we move the instance index from the subscript to the superscript, since we now use subscript for random seeds, and instance index can be omitted in most of our derivations.

The expected squared loss \mathcal{L} of model s on in-

stance i can then be written as

$$\mathcal{L}^{s,i} = \mathbb{E}_{P,F,E}[(1 - c_{P,F,E}^{s,i})^2]$$
 (33)

Since we will analyze this term at a datapoint level, we drop the subscript s and i to keep the notation uncluttered. By the standard bias variance decomposition and total variance decomposition, we decompose the loss \mathcal{L} into four terms:

$$\mathcal{L} = Bias^2 + PretVar$$
+ FineVar + CkptVar. (34)

We will walk through the meaning and definition of these four terms one by one. Bias² captures how bad is the average prediction, defined as

Bias² =
$$(1 - \mathbb{E}_{P.F.E}[c_{P.F.E}])^2$$
. (35)

PretVar captures the variance brought by randomness in pretraining, and is defined as

$$PretVar = Var_P[\mathbb{E}_{F,E}[c_{P,F,E}]].$$
 (36)

Similarly, we define the variance brought by randomness in finetuning FineVar

FineVar =
$$\mathbb{E}_P[Var_F[\mathbb{E}_E[c_{P,F,E}]]],$$
 (37)

and that by fluctuations across checkpoints e

$$CkptVar = \mathbb{E}_{P,F}[Var_E[c_{P,F,E}]]. \tag{38}$$

H.2 Unbiased Estimation

We first describe the data on which we apply our estimator. Suppose we pretrain \mathcal{P} models with different random seeds, for each of the \mathcal{P} pretrained models we finetune with \mathcal{F} different random seeds, and we evaluate at \mathcal{E} different checkpoints. Then $\forall j \in [\mathcal{P}], k \in [\mathcal{F}], l \in [\mathcal{E}]^{-13}$, we observe P_j, F_{jk}, E_{jkl} , and $c_{P_j, F_{jk}, E_{jkl}}$, where each observed P, F and E are i.i.d. distributed. Our goal is to estimate from c the four quantities described in the previous section.

H.2.1 Estimating CkptVar

It is straightforward to estimate CkptVar. The estimator CkptVar defined below is unbiased:

$$\operatorname{CkptVar} := \frac{1}{\mathcal{PF}} \sum_{j \in [\mathcal{P}], k \in [\mathcal{F}]} \hat{Var}_{E}^{P_{j}, F_{jk}}, \quad (39)$$

where

$$\hat{Var}_{E}^{P_{j},F_{jk}} := \frac{1}{\mathcal{E} - 1} \sum_{l \in \mathcal{E}} (c_{P_{j},F_{jk},E_{jkl}} - \bar{c}_{P_{j},F_{jk}})^{2},$$
(40)

and

$$\bar{c}_{P_j, F_{jk}} := \frac{1}{\mathcal{E}} \sum_{l \in [\mathcal{E}]} c_{P_j, F_{jk}, E_{jkl}}.$$
 (41)

CkptVar is unbiased, since $\hat{Var}_E^{P_j,F_{jk}}$ is an unbiased estimation of variance of c with fixed P_j and F_{jk} , and randomness E, i.e.

$$\mathbb{E}_{E_{ij(\cdot)}}[\hat{Var}_{E}^{P_{j},F_{jk}}] = Var_{E}[c_{P,F,E}|P = P_{j}, F = F_{jk}].$$
(42)

Therefore, $\forall j \in [\mathcal{P}], k \in [\mathcal{F}]$, we have

$$\mathbb{E}_{P_j, F_{jk}}[\hat{Var}_E^{P_j, F_{jk}}] = \text{CkptVar}, \quad (43)$$

and hence by linearity of expectation

$$\mathbb{E}_{P_{(\cdot)}, F_{(\cdot)}, E_{(\cdot)}}[\operatorname{CkptVar}] = \operatorname{CkptVar}.$$
 (44)

H.2.2 Estimating FineVar

As before, by linearity of expectation, we can declare victory if we can develop an unbiased estimator for the following quantity and then average across P_i :

$$Var_F[\mathbb{E}_E[c_{P,F,E}]|P=P_i],\tag{45}$$

which verbally means "variance across different finetuning seeds of the mean of c over different

$$^{13}[L] := \{l : l \in N, l \in [0, L-1]\}$$

checkpoints E, conditioned on the pretraining seed P_i ."

Since P_j is fixed for this estimator, we drop the subscripts P to keep notation uncluttered. Therefore, we want to estimate

$$Var_F := Var_F[\mathbb{E}_E[c_{F,E}]] \tag{46}$$

A naive solution is to take first take the mean \bar{c}_{F_k} of c for each F_k , i.e.

$$\bar{c}_{F_{jk}} := \frac{1}{\mathcal{E}} \sum_{l \in [\mathcal{E}]} c_{F_k, E_{kl}}, \tag{47}$$

and then calculate the sample variance \tilde{Var}_F of \bar{c} with respect to F:

$$\tilde{Var}_F := \frac{1}{\mathcal{F} - 1} \sum_{k \in [\mathcal{F}]} (\bar{c}_{F_k} - \bar{c})^2, \tag{48}$$

where

$$\bar{c} := \frac{1}{\mathcal{F}} \sum_{k \in [\mathcal{F}]} \bar{c}_{F_k} \tag{49}$$

However, this would create an upward bias: the empirical mean $\bar{c}_{F_{jk}}$ is a noisy estimate of the population mean $\mathbb{E}_E[c_{F_{jk},E}]$, and hence increases let \tilde{Var}_F over-estimate the variance. Imagine a scenario where Var_F is in fact 0; however, since $\bar{c}_{F_{jk}}$ is a noisy estimate, \tilde{Var}_F will sometimes be positive but never below 0. As a result, $\mathbb{E}[\tilde{Var}_F] > 0$, which is a biased estimator.

We introduce the following general theorem to correct this bias.

Theorem 2 Suppose $\mathcal{D}_k, k \in [\mathcal{F}]$ are independently sampled from the same distribution Ξ , which is a distribution of distributions; $\hat{\mu}_k$ is an unbiased estimator of $\mathbb{E}_{X \in \mathcal{D}_k}[X]$, and $\hat{\phi}_k$ to be an unbiased estimator of the variance of $\hat{\mu}_k$, then

$$\hat{Var}_F = \frac{1}{\mathcal{F} - 1} \sum_{k \in [\mathcal{F}]} (\hat{\mu}_k - \hat{\mu})^2$$

$$- \frac{1}{\mathcal{F}} \sum_{k \in [\mathcal{F}]} \hat{\phi}_k$$
(50)

is an unbiased estimator for

$$V = Var_{\mathcal{D} \sim \Xi}[\mathbb{E}_{X \sim \mathcal{D}}[X]], \tag{51}$$

where

$$\hat{\mu} := \frac{1}{\mathcal{F}} \sum_{k \in [\mathcal{F}]} \hat{\mu}_k \tag{52}$$

In this estimator, the first term "pretends" that $\hat{\mu}$. are perfect estimator for the population mean and calculate the variance, while the second term corrects for the fact that the empirical mean estimation is not perfect. Notice the theorem only requires that $\hat{\mu}$ and $\hat{\phi}$ are unbiased, and is agnostic to the actual computation procedure by these estimators.

Proof We define the population mean of \mathcal{D}_k to be μ_k , i.e.

$$\mu_k := \mathbb{E}_{X \sim \mathcal{D}_k}[X], \tag{53}$$

and the population mean of μ_k across randomness in \mathcal{D} to be μ , i.e.

$$\mu := \mathbb{E}_{\mathcal{D} \sim \Xi}[\mathbb{E}_{X \sim \mathcal{D}}[X]] \tag{54}$$

We look at the first term of the estimator in equation 50:

$$\frac{1}{\mathcal{F} - 1} \mathbb{E} \left[\sum_{k \in [\mathcal{F}]} (\hat{\mu}_k - \hat{\mu})^2 \right]$$

$$= \frac{1}{\mathcal{F} - 1} \mathbb{E} \left[\sum_{k \in [\mathcal{F}]} ((\hat{\mu}_k - \mu_k) - (\hat{\mu} - \mu)) + (\mu_k - \mu)^2 \right]$$

$$= \frac{1}{\mathcal{F} - 1} \mathbb{E} \left[\sum_{k \in [\mathcal{F}]} [(\hat{\mu}_k - \mu_k)^2 + (\hat{\mu} - \mu))^2 + (\mu_k - \mu)^2 - 2(\hat{\mu}_k - \mu_k)(\hat{\mu} - \mu) - 2(\hat{\mu}_k - \mu_k)(\hat{\mu} - \mu) \right]$$

There are 5 summands within $\sum_{k \in [\mathcal{F}]}$, and we look at them one by one:

$$\mathbb{E}\left[\sum_{k\in[\mathcal{F}]}(\hat{\mu}_k - \mu_k)^2\right] = \mathbb{E}\left[\sum_{k\in[\mathcal{F}]}\hat{\phi}_k\right],\tag{56}$$

$$\mathbb{E}[(\hat{\mu} - \mu)^2] = \mathbb{E}[(\mu - \frac{1}{\mathcal{F}} \sum_{k \in [\mathcal{F}]} \mu_k)$$

$$+ \frac{1}{\mathcal{F}} \sum_{k \in [\mathcal{F}]} (\mu_k - \hat{\mu}_k))^2]$$

$$= \frac{1}{\mathcal{F}} V + \frac{1}{\mathcal{F}^2} \sum_{k \in [\mathcal{F}]} \mathbb{E}[\hat{\phi}_k]$$
(57)

$$\mathbb{E}\left[\sum_{k\in[\mathcal{F}]}(\mu_k - \mu)^2\right] = \mathcal{F}V \tag{58}$$

$$\mathbb{E}[-2\sum_{k\in[\mathcal{F}]}(\hat{\mu}_k - \mu_k)(\hat{\mu} - \mu)]$$

$$= -\frac{2}{\mathcal{F}}\mathbb{E}[\sum_{k\in[\mathcal{F}]}\hat{\phi}_k].$$
(59)

$$\mathbb{E}[-2\sum_{k\in[\mathcal{F}]}(\hat{\mu}_k - \mu_k)(\hat{\mu} - \mu)]$$

$$= -2V. \tag{60}$$

Putting these five terms together, we continue calculating Equation 55:

$$\frac{1}{\mathcal{F} - 1} \mathbb{E} \left[\sum_{k \in [\mathcal{F}]} (\hat{\mu}_k - \hat{\mu})^2 \right] \tag{61}$$

$$= \frac{1}{\mathcal{F} - 1} \mathbb{E} \left[\sum_{k \in [\mathcal{F}]} \hat{\phi}_k \right]$$

$$+ \mathcal{F} \left(\frac{1}{\mathcal{F}} V + \frac{1}{\mathcal{F}^2} \sum_{k \in [\mathcal{F}]} \mathbb{E} \left[\hat{\phi}_k \right] \right)$$

$$+ \mathcal{F} V$$

$$- \frac{2}{\mathcal{F}} \sum_{k \in [\mathcal{F}]} \hat{\phi}_k$$

$$- 2V \right]$$

$$= V + \mathbb{E} \left[\frac{1}{\mathcal{F}} \sum_{k} \hat{\phi}_k \right]$$

Then from Equation 50, we can tell that $\hat{Var_F}$ is unbiased. \Box

Now we come back to the topic of developing an unbiased estimator for Var_F as defined in Equation 46. To utilize Theorem 2, we need two components:

- An unbiased estimator $\hat{\mu}_{F_k}$ for $\mathbb{E}_E[c_{F,E}|F=F_k]$
- An unbiased estimator $\hat{\phi}_{F_k}$ for the variance of $\hat{\mu}_{F_k}$, i.e. $Var_{E_{(\cdot)}}(\hat{\mu}_{F_k})$

 $ar{c}_{F_k}$ is an unbiased estimator for $\mathbb{E}_E[c_{F,E}|F=F_k]$, and its variance $Var_E[ar{c}_F|F=F_k]$ is

$$Var_{E(\cdot)}(\bar{c}_{F_k}) = \frac{1}{\mathcal{E}} Var_E[c_{F,E}|F = F_k]. \quad (62)$$

Therefore, to develop an unbiased estimator for $Var_E(\bar{c}_{F_{jk}})$, it suffices to have an unbiased estimate of $Var_E[c_{F,E}|F=F_k]$. We define

$$\hat{\phi}_{F_k} := \frac{1}{\mathcal{E}(\mathcal{F} - 1)} \sum_{k \in [\mathcal{L}]} (c_{F_k, E_{kl}} - \bar{c}_{F_k})^2, \quad (63)$$

and we can plug in $\hat{\phi}_{F_k}$ and $\hat{\mu}_{F_k} = \bar{c}_{F_k}$ into Theorem 2 as an unbiased estimator to obtain an unbiased estimator for $Var_F[\mathbb{E}_E[c_{P,F,E}]|P=P_j]$, and we average the estimation for each P_j to obtain an unbiased estimate.

H.2.3 Estimating PretVar

We next estimate $Var_P[\mathbb{E}_{F,E}[c_{P,F,E}]]$ We can still apply the idea from Theorem 2, which requires

- An unbiased estimator $\hat{\mu}_{P_j}$ for $\mathbb{E}_{F,E}[c_{P,F,E}|P=P_j]$
- An unbiased estimator $\hat{\phi}_{P_j}$ for the variance of $\hat{\mu}_{P_j}$, i.e. $Var_{F,E}[\hat{\mu}_{P_j}]$.

Again, the first is easy to obtain: $\hat{\mu}_{P_j} = \bar{c}_{P_j}$ is an unbiased estimator for $\mathbb{E}_{F,E}[c_{P,F,E}|P=P_j]$, where

$$\bar{c}_{P_j} := \frac{1}{\mathcal{F}\mathcal{E}} \sum_{k \in [\mathcal{F}], l \in [\mathcal{E}]} c_{P_j, F_{jk}, E_{jkl}}$$
(64)

However, we cannot straightforwardly estimate $Var_{F,E}[\hat{\mu}_{P_j}]$ as before, since samples $c_{P_j,F_{jk},E_{jkl}}$ are no longer independent. We need to use Equation 57 to develop an unbiased estimator (the LHS is exactly what we want!), i.e.

$$Var_{F,E}(\bar{c}_{P_j}) = \frac{1}{\mathcal{F}} Var_F(\mathbb{E}_E[c_{P,F,E}]|P = P_j)$$
(65)

$$+\frac{1}{\mathcal{F}^2}\sum_{k\in[\mathcal{F}]}Var_E(\bar{c}_{P_j,F_{jk}}),$$

and we already know how to estimate these two summands from the previous discussion on estimating FineVar.

H.2.4 Estimating $Bias^2$

It is easy to see that the following \hat{L} is an unbiased estimator for the loss \mathcal{L} .

$$\hat{\mathcal{L}} := \frac{1}{\mathcal{PFE}} \sum_{j \in [\mathcal{P}], k \in [\mathcal{F}], l \in [\mathcal{E}]} (1 - c_{P_j, F_{jk}, E_{jkl}})^2,$$
(66)

and

$$\mathbb{E}[\hat{\mathcal{L}}] = \mathcal{L}.\tag{67}$$

By linearity of expectation and loss decomposition in Equation 34,

$$\hat{\text{Bias}^2} := \hat{\mathcal{L}} - \text{PretVar}$$

$$- \hat{\text{FineVar}} - \hat{\text{CkptVar}}$$
(68)

is an unbiased estimator of Bias².

Notice that the naïve estimator that calculates the expected bias and then squares it estimates $(\mathbb{E}[\mathrm{Bias}])^2$ instead of $\mathbb{E}[\mathrm{Bias}^2]$.

H.3 Generalization

We can generalize this estimation strategy to decompose variance into arbitrary number of randomness. In general, we want to estimate some quantity of the following form

$$\mathbb{E}_{r_1...,r_{n-1}}[Var_{r_n}[\mathbb{E}_{r_{n+1}...r_N}[c_{r_1...,c_N}]]], \quad (69)$$

from the data that has an hierarchical tree structure of randomness.

For the goal of developing an unbiased estimator, we can get rid of the outer expectation $r_1 cdots r_{n-1}$ easily by linearity of expectation: simply estimate the Variance conditioned on $r_{1...n-1}$ and average them together, as discussed in Section H.2.1.

To estimate

$$Var_{r_n}[\mathbb{E}_{r_{n+1}\dots r_N}[c_{r_1\dots c_N}]], \tag{70}$$

we make use of Theorem 2, which requires

- an unbiased estimator $\hat{\mu}_{r_{n+1}}$ for the quantity $\mathbb{E}_{r_{n+1}...r_N}[c_{r_1...,c_N}]$, which we can straightforwardly obtain by average the examples that has the same random variables $r_{1...n}$ (e.g. \bar{c}_{P_i})
- an unbiased estimator for the variance of $\hat{\mu}_{r_{n+1}}$. If N=n+1, we can directly compute the sample variance of the c as our estimate (e.g. in Equation 63). Otherwise, we use Equation 57 to decompose the desired quantities into two, and estimate them recursively by applying Theorem 2 and Equation 57.

For readability we wrote the proof with the assumption that, in the tree of randomness, the number of branches for each node at the same depth is the same. However, our proof does not make use of this assumption and can be applied to a general tree structure of randomness as long as the the number of children is larger or equal to 2 for each non-terminal node.

I Variance Conditioned on Bias

Since lower bias usually implies lower variance, to tease out the latent effect, we estimate the variance "given a fixed level of bias Bias^2 of $b^2 \in [0,1]$ ", i.e.

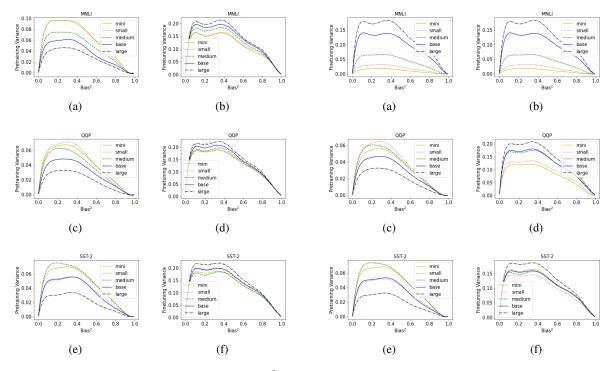


Figure 9: The variance curve conditioned on Bias² for in-domain development set of MNLI, QQP and SST-2. Each curve represents a model size. Left for pretraining variance and right for finetuning variance.

Figure 10: The same figure as 9, except for using the squared loss function $\mathcal{L} = (1-p)^2$, where p is the probability assigned to the correct label, instead of 0/1 loss.

$$\operatorname{PretVar}(b^2) := \mathbb{E}_i[\operatorname{PretVar}_i|\operatorname{Bias}^2_i = b^2]$$
 (71)

We estimate $\operatorname{PretVar}^s(b^2)$ and $\operatorname{FineVar}^s(b^2)$ using gaussian process and plot them against b^2 in Figure 9 for MNLI, QQP, and SST-2. We find that larger models usually have larger finetuning variance across all levels of biases (except for MEDIUM and MINI on SST-2), and BASE model always has larger pretraining variance than LARGE .

We also experimented with the squared loss:

$$\mathcal{L}_i = (1 - p_i)^2,\tag{72}$$

where p_i is the probability the assigned to the correct label for instance i. We plot the same curve in Figure 10 and observe the same trend.

J Example Decaying Instances

We manually examined the group of instances where $^{\rm MINI}_{\rm LARGE}\Delta\hat{\rm A}{\rm cc}_i \le -0.9$ in Table 3. In other words, MINI is almost always correct on these instances, while LARGE is almost always wrong. For each instance in this group, we manually categorize it into one of the four categories: 1) Correct, if the label is correct, 2) Fine, if the label might

be controversial but we could see a reason why this label is reasonable, 3) Wrong, if the label is wrong, and 4) Unsure, if we are unsure how to label this instance. As a control, we also examined the remaining fraction of the dataset. Each time we annotate an instance, with 50% probability it is sampled from the decaying fraction or the remaining fraction, and we do not know which group it comes from.

We show below all the annotated instances from this decaying fraction and their categories for MNLI (Section J.1), QQP, and SST-2(Section J.3).

J.1 MNLI

MNLI is the abbreviation of Multi-Genre Natural Language Inference (Williams et al. (2020)). In this task, given a premise and a hypothesis, the model needs to classify whether the premise entails/contradicts the hypothesis, or otherwise. The instances can be seen below.

Premise: and that you're very much right but the jury may or may not see it that way so you get a little anticipate you know anxious there and go well you know

Hypothesis: Jury's operate without the benefit of

an education in law. **Label**: Neutral **Category**: Correct

Premise: In fiscal year 2000, it reported estimated improper Medicare Fee-for-Service payments of \$11.

Hypothesis: The payments were improper.

Label: Entailment **Category**: Fine

Premise: is that what you ended up going

into

Hypothesis: So that must be what you chose to

do?

Label: Entailment **Category**: Correct

Premise: INTEREST RATE - The price charged per unit of money borrowed per year, or other unit of time, usually expressed as a percentage.

Hypothesis: Interest rate is defined as the total amount of money borrowed.

Label: Entailment **Category**: Wrong

Premise: The analyses comply with the informational requirements of the sections including the classes of small entities subject to the rule and alternatives considered to reduce the burden on the small entities.

Hypothesis: The rules place a high burden on the activities of small entities.

Label : Contradiction
Category : Correct

Premise: Isn't a woman's body her most personal property?

Hypothesis: Women's bodies belong to themselves, they should decide what to do with it.

Label: Neutral **Category**: Unsure

Premise: The Standard, published a few days before Deng's death, covers similar territory. **Hypothesis**: The Weshington Post severe similar

Hypothesis: The Washington Post covers similar

territory. **Label**: Neutral

Category : Correct

Premise: Shoot only the ones that face us,

Jon had told Adrin.

Hypothesis: Jon told Adrin and the others to only

shoot the ones that face us.

Label: Entailment **Category**: Wrong

Premise: But if you take it seriously, the anti-abortion position is definitive by definition. **Hypothesis**: If you decide to be serious about

supporting anti-abortion, it's a very run of the mill

belief to hold. **Label**: Neutral **Category**: Unsure

Premise: yeah well that's the other thing you know they talk about women leaving the home and going out to work well still taking care of the children is a very important job and and someone's got to do it and be able to do it right and

Hypothesis: It is not acceptable for anybody to refuse work in order to take care of children.

Label : Contradiction **Category** : Correct

Premise: The researchers found expected stresses like the loss of a check in the mail and the illness of loved ones.

Hypothesis: The stresses affected people much diffferently than the researchers expected.

Label: Contradiction **Category**: Correct

Premise: so you know it's something we

we have tried to help but yeah

Hypothesis: We did what we could to help.

Label: Entailment **Category**: Correct

Premise: Czarek was welcomed enthusiastically, even though the poultry brotherhood was paying a lot of sudden attention to the newcomers - a strong group of young and talented managers from an egzemo-exotic chicken farm in Fodder Band nearby Podunkowice.

Hypothesis: Czarek was welcomed into the group by the farmers.

Label : Entailment Category : Correct

Premise: 'I don't suppose you could forget

I ever said that?'

 $\boldsymbol{Hypothesis}: I \ hope \ that \ you \ can \ remember \ that$

forever.

Label : Contradiction **Category** : Wrong

Premise: Oh, my friend, have I not said to

you all along that I have no proofs.

Hypothesis: I told you from the start that I had no

evidence.

Label: Entailment **Category**: Correct

Premise: I should put it this way.

Hypothesis: I should phrase it differently.

Label: Entailment **Category**: Correct

Premise: An organization's activities, core processes, and resources must be aligned to support its mission and help it achieve its goals.

Hypothesis: An organization is successful if its

activities, resources, and goals align.

Label: Entailment **Category**: Fine

Premise: A more unusual dish is azure, a kind of sweet porridge made with cereals, nuts, and fruit sprinkled with rosewater.

Hypothesis: Azure is a common and delicious

food made with cereals, nuts and fruit.

Label: Entailment **Category**: Wrong

Premise: once you have something and it's like i was watching this program on TV yesterday in nineteen seventy six NASA came up with Three D graphics right

Hypothesis: I was watching a program about

gardening.

Label : Contradiction **Category** : Correct

Premise:, First-Class Mail used by households to pay their bills) and the household bill mail (i.e.

Hypothesis: Second-Class Mail used by house-

holds to pay their bills **Label**: Contradiction **Category**: Unsure

Premise: Rightly or wrongly, America is seen as globalization's prime mover and head cheerleader and will be blamed for its excesses until we start paying official attention to them.

Hypothesis: America's role in the globalization movement is important whether we agree with it or

Label: Entailment **Category**: Correct

Premise: After being diagnosed with cancer, Carrey's Kaufman decides to do a show at Carnegie Hall.

Hypothesis: Carrey's Kaufman is only diagnosed with cancer after doing a show at Carnegie Hall.

Label: Contradiction **Category**: Correct

Premise: Several pro-life Dems are mounting serious campaigns at the state level, often against pro-choice Republicans.

Hypothesis: Serious campaigns are being run by

a few pro-life Democrats.

Label: Entailment **Category**: Correct

Premise: On the northwestern Alpine frontier, a new state had appeared on the scene, destined to lead the movement to a united Italy.

Hypothesis: The unite Italy movement was waiting for a leader.

Label: Neutral
Category: Fine

Premise: well we bought this with credit too well we found it with a clearance uh down in Memphis i guess and uh

Hypothesis: We bought non-sale items in

Memphis on credit. **Label**: Contradiction **Category**: Correct

Premise: He slowed.

Hypothesis: He stopped moving so quickly.

Label: Entailment **Category**: Correct

Premise: As legal scholar Randall Kennedy wrote in his book Race, Crime, and the Law, Even if race is only one of several factors behind a decision, tolerating it at all means tolerating it as

potentially the decisive factor.

Hypothesis: Race is one of several factors in

some judicial decisions **Label** : Entailment **Category** : Correct

Premise: Although all four categories of emissions are down substantially, they only achieve 50-75% of the proposed cap by 2007 (shown as the dotted horizontal line in each of the above figures).

Hypothesis: All of the emission categories experienced a downturn except for one.

Label: Contradiction **Category**: Correct

Premise: He sat up, trying to free himself. **Hypothesis**: He was trying to take a nap.

Label: Contradiction **Category**: Correct

Premise: Impossible.

Hypothesis: Cannot be done.

Label: Entailment **Category**: Correct

Premise: But, as the last problem I'll outline suggests, neither of the previous two objections matters.

Hypothesis: I will not continue to outline any more problems.

Label: Entailment
Category: Correct

Premise: As the Tokugawa shoguns had feared, this opening of the floodgates of Western culture after such prolonged isolation had a traumatic effect on Japanese society.

Hypothesis: The Tokugawa shoguns had feared that, because they understood the Japanese society very well.

Label: Neutral Category: Fine

Premise: In the ancestral environment a man would be likely to have more offspring if he got his pick of the most fertile-seeming women.

Hypothesis: Only a man who stayed with one female spread his genes most efficiently.

Label : Contradiction **Category** : Fine

Premise: Tommy was suddenly galvanized

into life.

Hypothesis: Tommy had been downcast for days.

Label: Neutral **Category**: Correct

Premise: Improved products and services Initiate actions and manage risks to develop new products and services within or outside the organization.

Hypothesis: Managed risks lead to new products

Label: Entailment **Category**: Fine

Premise : Coast Guard rules establishing

bridgeopening schedules).

Hypothesis: The Coast Guard is in charge of

opening bridges. **Label**: Entailment **Category**: Correct

Premise: The anthropologist Napoleon Chagnon has shown that Yanomamo men who have killed other men have more wives and more offspring than average guys.

Hypothesis: Yanomamo men who kill other men have better chances at getting more wives.

Label: Entailment **Category**: Fine

Premise: The Varanasi Hindu University has an Art Museum with a superb collection of 16th-century Mughal miniatures, considered superior to the national collection in Delhi.

Hypothesis: The Varanasi Hindu University has an art museum on its campus which may be superior objectively to the national collection in Delhi.

Label: Entailment **Category**: Correct

Premise: Part of the reason for the difference in pieces per possible delivery may be due to the fact that five percent of possible residential deliveries are businesses, and it is thought, but not known, that a lesser percentage of possible deliveries on rural routes are businesses.

Hypothesis: We all know that the reason for a lesser percentage of possible deliveries on rural routes being businesses, is because of the fact that

people prefer living in cities rather than rural areas.

Label : Neutral
Category : Correct

Premise: right oh they've really done uh good job of keeping everybody informed of what's going on sometimes i've wondered if it wasn't almost more than we needed to know

Hypothesis: I don't think I have shared enough information with everyone.

Label: Contradiction
Category: Correct

Premise: To reach any of the three Carbet falls, you must continue walking after the roads come to an end for 20 minutes, 30 minutes, or two hours respectively.

Hypothesis: There are three routes to the three Carbet falls, each a different length and all continue after the road seemingly ends.

Label: Entailment **Category**: Correct

Premise: But when the cushion is spent in a year or two, or when the next recession arrives, the disintermediating voters will find themselves playing the roles of budget analysts and tax wonks.

Hypothesis: The cushion will likely be spent in under two years.

Label : Entailment **Category** : Correct

Premise: But, Slate protests, it was [Gates'] byline that appeared on the cover.

Hypothesis: Slate was one hundred percent positive it was Gates' byline on the cover.

Label : Neutral
Category : Correct

Premise: But it's for us to get busy and do something."

Hypothesis: "We don't do much, so maybe this would be good for us to bond and be together for the first time in a while.".

Label: Neutral **Category**: Fine

Premise: Pearl Jam detractors still can't stand singer Eddie They say he's unbearably self-important and limits the group's appeal by refusing to sell out and make videos.

Hypothesis: A lot of people consider Eddie to be

a bad singer. **Label**: Neutral **Category**: Correct

Premise: it's the very same type of paint

and everything

Hypothesis: It's the same paint formula, it's

great!

Label: Entailment **Category**: Fine

Premise: Exhibit 3 presents total national emissions of NOx and SO2 from all sectors, including power.

Hypothesis: In Exhibit 3 there are the total regional emissions od NOx and SO2 from all sectors.

Label: Entailment **Category**: Correct

Premise: uh-huh and is it true i mean is it

um

Hypothesis: It's true.
Label: Entailment
Category: Wrong

Premise: When a GAGAS attestation engagement is the basis for an auditor's subsequent report under the AICPA standards, it would be advantageous to users of the subsequent report for the auditor's report to include the information on compliance with laws and regulations and internal control that is required by GAGAS but not required by AICPA standards.

Hypothesis: The report is required by GAGAS but not AICPA.

Label: Entailment
Category: Correct

Premise: i'm on i'm in the Plano school system and living in Richardson and there is a real dichotomy in terms of educational and economic background of the kids that are going to be attending this school

Hypothesis: The Plano school system only has children with poor intelligence.

Label: Contradiction **Category**: Correct

J.2 QQP

QQP is the abbreviation of Quora Question Pairs¹⁴. Given two questions, the model needs to tell whether they have the same meaning (i.e. Paraphrase/Non-paraphrase).

Question 1: Which universities for MS in CS should I apply to?

Question 2: Which universities should I apply to for an MS in CS?

Label : Paraphrase Category : Correct

Question 1: What should I do to make life worth living?

Question 2: What makes life worth living?

Label: Paraphrase **Category**: Fine

Question 1: Why did Quora remove my question?

Question 2: Why does Quora remove questions?

Label: Paraphrase **Category**: Correct

Question 1: How do I get thousands of followers on Instagram?

Question 2: How can I get free 10k real Instagram followers fast?

Label: Paraphrase Category: Fine

Question 1: What is the basic knowledge of computer science engineers?

Question 2: What is basic syallbus of computer science engineering?

Label: Non-paraphrase

Category: Fine

Question 1: How many mosquito bites does it take to kill a human being?

Question 2: How many times can a single mosquito bite a human within 8 hours?

Label: Non-paraphrase **Category**: Correct

Question 1: How does it feel to become attractive from unattractive?

Question 2: What does it feel like to go from physically unattractive to physically attractive?

Label: Paraphrase **Category**: Correct

Question 1: Who is answering the questions asked on Quora?

Question 2: Who can answer the questions asked on Ouora?

Label : Paraphrase Category : Correct

Question 1: What machine learning theory do I need to know in order to be a successful machine learning practitioner?

Question 2: What do I need to know to learn machine learning?

Label: Paraphrase **Category**: Wrong

Question 1: If you could go back in time and change one thing, what would it be and why?

Question 2: If you could go back in time and do one thing, what would it be?

Label: Paraphrase **Category**: Correct

Question 1: Will there be a civil war if Trump doesn't become president?

Question 2: Will there be a second civil war if Trump becomes president?

Label: Paraphrase **Category**: Correct

Question 1: Do Quora contributors get paid?

Question 2: How do contributors get paid by Ouora?

Label: Paraphrase **Category**: Correct

Question 1: Did India meet Abdul Kalam's 2020 vision so far?

Question 2: How far do you think India has reached on President APJ Kalam's vision in the book India 2020?

Label : Non-paraphrase **Category** : Correct

Question 1: How do I stop my dog from whining after getting spayed?

Question 2: How do I stop my dog from whining?

Label: Paraphrase

¹⁴https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs

Category: Wrong

Question 1: What difference are exactly between Euclidean space and non Euclidean space?

Question 2: What is the difference between Euclidean and non-Euclidean?

Label : Non-paraphrase **Category** : Wrong

Question 1: Why doesn't Hillary Clinton win the White House if she won the popular vote?

Question 2: How did Hillary Clinton win the popular vote but Donald Trump win the election?

Label : Paraphrase **Category** : **Correct**

Question 1: How is public breastfeeding seen where you live?

Question 2: How is breastfeeding in public seen in your country?

Label: Paraphrase
Category: Correct

Question 1: What are some ways to change your Netflix password?

Question 2: How do you change your Netflix password and email?

Label: Paraphrase
Category: Fine

Question 1: What do you think, is your best answer on Quora?

Question 2: What is your best answer on Quora?

Label : Paraphrase **Category** : Correct

Question 1: How can I travel to Mexico without a passport?

Question 2: Can I travel to Mexico without a passport?

Label: Paraphrase **Category**: Correct

Question 1: How do modern Congolese people view Mobutu in retrospect?

Question 2: How do Congolese currently view

Mobutu Sese Seko? **Label**: Non-paraphrase **Category**: Correct

Question 1: How is Tanmay Bhat losing weight?

Question 2: Tanmay Bhat: How did you manage to reduce your fat?

Label: Non-paraphrase **Category**: Unsure

Question 1: Is Xiaomi a brand to trust (comparing it with brands like Samsung and HTC)? What is better: Xiaomi MI3 or HTC Desire 816?

Question 2: Is xiaomi a trusted brand?

Label: Non-paraphrase **Category**: Correct

Question 1: Why did Buddhism spread in East Asia and not in its native land India?

Question 2: How was Buddhism spread in Asia?

Label: Non-paraphrase **Category**: Correct

Question 1: Can I become a multi billionaire betting on horses?

Question 2: How much money can I make betting on horses? A month? Can I make 20,000 a month?

Label: Paraphrase **Category**: Fine

Question 1: What is a diet for gaining weight?

Question 2: What is a way to gain weight?

Label : Non-paraphrase **Category** : Correct

Question 1: How do I use Instagram on my computer?

Question 2: How can I get Instagram on my computer?

Label: Paraphrase **Category**: Fine

Question 1: What is the legal basis of a "you break it, you buy it" policy?

Question 2: Is a "you break it you buy it" policy actually legal?

Label: Paraphrase **Category**: Correct

Question 1: How I should fix my computer while it is showing no boot device found?

Question 2: How do I fix the "Boot device not

found" problem? **Label**: Paraphrase **Category**: Correct

Question 1: What innovative name can I use for an interior designing firm?

Question 2: What can i name my interior designing firm?

Label: Paraphrase Category: Correct

Question 1: What would it realistically cost to go to Tomorrowland?

Question 2: How much is a ticket to Tomorrowland?

Label: Non-paraphrase

Category: Fine

Question 1: Is there a gender pay gap? If so why?

Question 2: Is the gender pay gap a myth?

Label: Paraphrase **Category**: Correct

Question 1: How can I get rid of a canker sore on the bottom of my tongue?

Question 2: How can I get rid or a canker sore on the tip of my tongue?

Label : Paraphrase
Category : Correct

Question 1: How can I sleep better and early in night?

Question 2: How can I sleep better at night?

Label: Paraphrase **Category**: Fine

Question 1: Why did DC Change Captain Marvel's name?

Question 2: Why did DC have to change Captain Marvel's name but Marvel didn't have to change Scarecrow's name?

Label: Paraphrase Category: Fine

Question 1: Should be there any difference between IIT and non IIT students in terms of placement package from a company if both of them are equally talented?

Question 2: Should be there any difference between IIT and non IIT students in terms of

placement package from a company if both of them are equally capable?

Label: Paraphrase **Category**: Correct

Question 1: What happened to The Joker after The end of The Dark Knight?

Question 2: What happens to the Joker at the end

of The Dark Knight (2008 movie)? **Label**: Non-paraphrase

Category: Wrong

Question 1: I love my wife more then anything. Why do I fantasize about her with other men?

Question 2: Why do I fantasize about other men having sex with my wife?

Label: Paraphrase **Category**: Fine

Question 1: What is your opinion on the new MacBook Pro Touch Bar?

Question 2: What do you think about the OLED touch bar on the new MacBook Pro?

Label: Paraphrase **Category**: Correct

Question 1: How do I get rid of my negative alter ego?

Question 2: How do you get rid of your negative alter ego?

Label: Paraphrase **Category**: Correct

Question 1: How can I get wifi driver for my hp laptop with windows 7 os?

Question 2: How can I get wifi driver for my laptop with windows 7 os?

Label: Paraphrase **Category**: Fine

Question 1: What's your attitude towards life?

Question 2: What should be your attitude towards life?

Label: Paraphrase **Category**: Wrong

Question 1: What books would I like if I loved A Song of Ice and Fire?

Question 2: Are there books which are similar to

A Song of Ice and Fire? **Label**: Paraphrase Category : Fine

Question 1: Why do Muslims think they will conquer the whole world?

Question 2: Do you think Muslims will take over the world?

Label: Non-paraphrase **Category**: Correct

Question 1: Is dark matter a sea of massive dark photons that ripple when galaxy clusters collide and wave in a double slit experiment?

Question 2: Does a superfluid dark matter which ripples when Galaxy clusters collide and waves in a double slit experiment relate GR and QM?

Label: Paraphrase **Category**: Correct

Question 1: What is Batman like?

Question 2: What is Batman's personality like?

Label: Non-paraphrase **Category**: Correct

J.3 SST-2

SST-2 is the abbreviation of Stanford Sentiment Treebank (Socher et al., 2013). In this task, the model needs to recognize whether the phrases or sentences reflect positive or negative sentiments.

Input: predictability is the only winner

Label: Negative **Category**: Correct

Input: abandon their scripts and go where

the moment takes them Label: Negative **Category**: Correct

Input: chases for an hour and then

Label: Positive **Category**: Unsure

Input: provide much more insight than the

inside column of a torn book jacket

Label: Negative Category: Unsure **Input**: a children 's party clown

Label: Negative Category: Fine

Input: perhaps even the slc high command found writer-director mitch davis 's wall of kitsch hard going.

Label: Negative **Category**: Correct

Input: own placid way Label: Negative

Category: Correct

Input: get on a board and, uh, shred

Label: Negative **Category**: Correct

Input: asks what truth can be discerned from non-firsthand experience, and specifically questions cinema 's capability for recording truth .

Label: Positive **Category**: Correct

Input: puts the dutiful efforts of more dis-

ciplined grade-grubbers

Label: Positive **Category**: Correct

Input: filter out the complexity

Label: Positive **Category**: Correct

Input: told what actually happened as if it

were the third ending of clue

Label: Negative **Category**: Correct

Input: is more in love with strangeness

than excellence. **Label**: Positive **Category**: Wrong

Input: i found myself howling more than

cringing

Label: Positive **Category**: Correct

Input: goldbacher draws on an elegant visual sense and a talent for easy, seductive pacing ... but she and writing partner laurence coriat do n't manage an equally assured narrative coinage

Label: Positive **Category**: Unsure

Input: for a thirteen-year-old 's book re-

port

Label: Negative **Category**: Correct

Input: a problem hollywood too long has

ignored

Label: Negative **Category**: Correct

Input: twisted sense
Label: Negative
Category: Correct

Input: a stab at soccer hooliganism

Label: Negative **Category**: Correct

Input: sinuously plotted

Label: Negative **Category**: Correct

Input: shiner can certainly go the distance

, but is n't world championship material

Label: Positive Category: Correct

Input: holding equilibrium up

Label: Negative **Category**: Wrong

Input: i am highly amused by the idea that we have come to a point in society where it has been deemed important enough to make a film in which someone has to be hired to portray richard dawson.

Label: Positive **Category**: Wrong

Input: waters
Label: Negative
Category: Wrong

Input: what might have emerged as hilarious lunacy in the hands of woody allen or

Label: Positive

Category: Correct

Input: of those airy cinematic bon bons whose aims – and by extension, accomplishments

seem deceptively slight on the surface

Label: Positive **Category**: Correct

Input: do n't blame eddie murphy but

Label: Negative **Category**: Correct

Input: melodramatic paranormal romance

Label: Negative **Category**: Correct

Input: could possibly be more contemptu-

ous of the single female population

Label: Negative **Category**: Correct

Input: cremaster 3 " should come with the

warning "for serious film buffs only!"

Label: Negative **Category**: Correct

Input: softheaded metaphysical claptrap

Label: Negative **Category**: Correct

Input : owed to benigni
Label : Negative
Category : Unsure

Input: to be a suspenseful horror movie or

a weepy melodrama **Label**: Positive **Category**: Correct

Input: genuinely unnerving.

Label: Positive **Category**: Correct

Input: gaping enough to pilot an entire

olympic swim team through

Label: Negative **Category**: Correct

Input: this is popcorn movie fun with equal doses of action, cheese, ham and cheek (as well as a serious debt to the road warrior), but it feels like

unrealized potential **Label**: Positive **Category**: Fine

Input: feeling like it was worth your seven bucks, even though it does turn out to be a bit of a

cheat in the end **Label**: Negative **Category**: Correct

Input : pull it back
Label : Negative
Category : Correct

Input:, this is more appetizing than a side

dish of asparagus . **Label** : Negative **Category** : Correct

Input : crime drama Label : Negative Category : Unsure

Input: like most movies about the pitfalls

of bad behavior **Label**: Negative **Category**: Fine

Input: befallen every other carmen before

her

Label: Positive **Category**: Unsure

Input: appeal to those without much interest in the elizabethans (as well as rank frustration from those in the know about rubbo 's dumbed-

down tactics) **Label** : Negative **Category** : Unsure

Input: about existential suffering

Label: Negative **Category**: Fine

Input : , if uneven ,
Label : Negative
Category : Unsure

Input: succumbs to sensationalism

Label: Positive **Category**: Wrong

Input: that turns me into that annoying specimen of humanity that i usually dread encoun-

tering the most **Label**: Negative **Category**: Fine

Input: at least a minimal appreciation

Label: Positive **Category**: Unsure

Input: underlines even the dullest tangents

Label : Negative **Category** : Correct

Input : heard before
Label : Positive
Category : Unsure

Input: i like my christmas movies with more elves and snow and less pimps and ho 's.

Label: Negative **Category**: Unsure

Input: can aspire but none can equal

Label: Negative **Category**: Unsure

Input: fathom
Label: Negative
Category: Unsure

Input: attempt to bring cohesion to pamela

's emotional roller coaster life

Label: Negative **Category**: Unsure

Input: movie version Label: Positive Category: Wrong

Input: of spontaneity in its execution and

a dearth of real poignancy

Label: Positive **Category**: Correct