
Predicting the Performance of Foundation Models via Agreement-on-the-Line

Aman Mehra*¹ Rahul Saxena*¹ Taeyoun Kim*¹ Christina Baek¹
Zico Kolter^{1,2} Aditi Raghunathan¹

Carnegie Mellon University¹, Bosch Center for AI²

{amanmehra, rsaxena2, taeyoun3, kbaek, zkolter, raditi}@andrew.cmu.edu

Abstract

1 Estimating out-of-distribution performance is critical to safely deploy machine
2 learning models. Recently, Baek et al. showed that the phenomenon “agreement-
3 on-the-line” can be a reliable method for predicting the OOD accuracy of models
4 in an ensemble consisting largely of CNNs trained from scratch. However, it is
5 now increasingly common to lightly fine-tune foundation models, and it is unclear
6 whether such fine-tuning is sufficient to produce enough diversity in model predic-
7 tions for such agreement-based methods to work properly. In this paper, we develop
8 methods for reliably applying agreement-on-the-line-based performance estimation
9 to fine-tuned foundation models. In particular, we first study the case of fine-tuning
10 a single foundation model, where we extensively study how different types of ran-
11 domness (linear head initialization, data shuffling, and data subsetting) contribute
12 to agreement-on-the-line of the resulting model sets. Somewhat surprisingly, we
13 find that it is possible to obtain strong agreement via random initialization of the
14 linear head alone. Next, we find how *multiple* foundation models, pretrained on dif-
15 ferent data sets but fine-tuned on the same task, also observe agreement-on-the-line.
16 Again rather surprisingly, the diversity of such models is not too disparate, and
17 they all lie on the same agreement line. In total, these methods enable reliable and
18 efficient estimation of OOD accuracy for fine-tuned foundation models, without
19 leveraging any labeled OOD data.

20 1 Introduction

21 Foundation models (FM) approaches, where one first pretrains a large model on open world data then
22 fine-tunes for a specific downstream task, have achieved state-of-the-art results on image classification
23 [27, 21, 38], text classification [6], question answering [8], and others. They are particularly noted
24 for their often strong performance on out-of-distribution (OOD) data, that may vary substantially
25 from the data used for fine-tuning (referred to as the in-distribution (ID) data) [5, 39]. Unfortunately,
26 a substantial practical problem arises in this OOD setting: in many cases, one does not have access to
27 labeled OOD data, and thus the field has explored other means for estimating OOD accuracy.

28 Interestingly, across a variety of distribution shift benchmarks, models often observe strong linear
29 correlation between the ID and OOD accuracies of models, a phenomenon dubbed “Accuracy-on-the-
30 line” (ACL) [25, 31, 32]. Recently, Baek et al. [2] empirically demonstrated that for ensembles of
31 deep network classifiers trained from scratch, the rates of ID and OOD agreement also show a strong
32 linear correlation with the same slope and bias. Baek et al. [2] used this to estimate the accuracies
33 of models in such ensembles, thus providing a simple method for estimating OOD accuracy via
34 unlabeled data alone. Thus, whenever the ID versus OOD accuracy is strongly linearly correlated,
35 one may estimate the linear OOD performance trend using agreement without ground truth labels.

36 Unfortunately, the AGL approach requires a *diverse collection* of classifiers over which to compute
37 agreement: classifiers must vary in their predictions. Baek et al. [2] achieve this by training various
38 models of different architectures from scratch. However, in the case of fine-tuned FMs, this diversity
39 is seemingly lacking: we often want to *lightly* fine-tune just a single base FM for a downstream
40 task, which even after multiple runs would seemingly lead to highly correlated downstream models,
41 making them unsuitable for AGL-based OOD performance estimation.

42 In this work, we develop methods for extending AGL performance estimation to FMs, thus enabling
43 practitioners to estimate the OOD performance of fine-tuned models without any labeled data. We
44 first investigate the ability to estimate performance using a *single* base FM. We present a detailed
45 empirical study of three potential sources of randomness during fine-tuning: 1) random linear head
46 initialization; 2) random orderings of the fine-tuning data; and 3) random i.i.d subsets of the fine-
47 tuning data. We find, somewhat surprisingly, that using random linear heads is able to reliably induce
48 AGL behavior for the resulting classifiers, with the result holding across multiple different FMs and
49 modalities (image classification and question answering a.k.a QA tasks). The result is a simple and
50 straightforward method for evaluating OOD performance for a fine-tuned FM, applicable to settings
51 where we only one want to fine-tune a single such base FM.

52 Second, we analyze the ability of the AGL-based method to predict OOD performance when using
53 *multiple* different pretrained FMs. Here we encounter a setting where the different base models
54 are pretrained on potentially entirely different data sets, using different architectures, and different
55 training regiments. We show, however, that this degree of diversity is *also* sufficient for producing
56 AGL behavior. Thus, for settings where multiple pretrained models exist, they can all be fine-tuned
57 for a given downstream task, and AGL can allow us to estimate their accuracies.

58 In total, our contributions are as follows:

- 59 1. We propose a state-of-the-art method for unsupervised accuracy estimation under distribution
60 shift when using large pretrained foundation models that are lightly fine-tuned for specific
61 tasks. Prior works have primarily dealt with models trained from scratch, and hence are not
62 directly applicable in this setting.
- 63 2. Our work leverages Agreement-on-the-line (AGL) [2] for OOD estimation, but extends it in
64 important ways to apply to finetuned foundation models. The key to making AGL work is
65 obtaining the right ensemble. In Baek et al. [2], multiple models were trained independently
66 from scratch, an unfeasible step for FMs. We show that creating an ensemble with randomly
67 initialized linear heads and then fine-tuning, also allows for AGL behavior, while other
68 similar forms of ensembling (such as data ordering or data subsetting) do not.
- 69 3. We also identify several interesting phenomena underlying AGL that go beyond previous
70 knowledge. Prior work Baek et al. [2] claimed that AGL does not hold for linear models.
71 However, we find the contrary when using pretrained CLIP features. Furthermore, other
72 prior work Miller et al. [25] suggests that the effective robustness (i.e. the linear fit between
73 ID and OOD accuracy) would change depending on the pretraining data. We find that this is
74 not the case for question answering with different pretrained FMs.

75 In total, this work substantially expands the set of problems and models for which AGL-based OOD
76 performance estimation is practical, and allows us to leverage much more powerful models for
77 settings where training models from scratch on tasks of interest is not feasible.

78 2 Background and related work

79 **OOD performance estimation of FMs.** Numerous tasks of interest boil down to mapping an input
80 $x \in \mathbb{X}$ to a discrete output $y \in \mathbb{Y}$. In particular, consider a base FM $B : \mathbb{X} \mapsto \mathbb{R}^d$ that we fine-tune to
81 get $f(B) : \mathbb{X} \mapsto \mathbb{Y}$. In this work, we consider a variety of foundation models: BERT [9], GPT2 [27],
82 GPT-Neo, OPT [41], Llama2 [36], and CLIP [28].

83 Given access to a labeled validation set from \mathcal{D}_{ID} and *unlabeled* samples from a different distribution
84 \mathcal{D}_{OOD} , our goal is to estimate performance on \mathcal{D}_{OOD} . We consider the standard performance metrics:
85 Accuracy $\ell_{0,1} : \mathbb{Y} \mapsto \mathbb{Y}$ for classification, and Macro-averaged F1 score $\ell_{F1} : \mathbb{Y} \mapsto \mathbb{Y}$ for QA.

86 There are a variety of proposed approaches for OOD performance estimation. One family of
87 approaches attempts to quantify the degree of distribution shift through data and/or model dependent

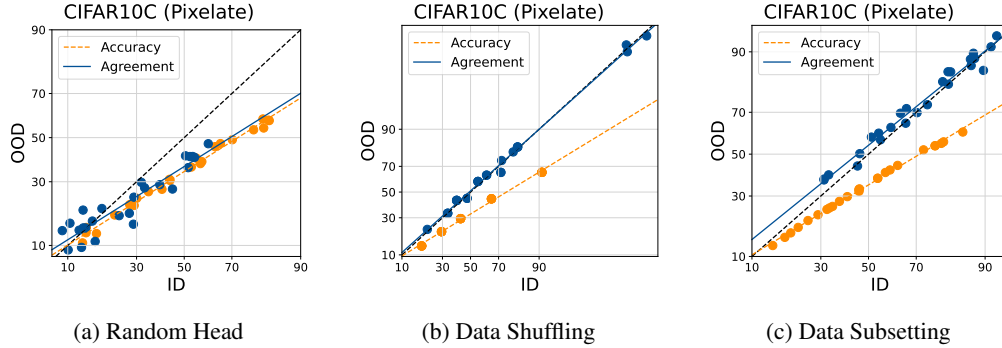


Figure 1: The ID vs OOD trends for accuracy and agreement on the CIFAR10C “Pixelate” shift for linear-probe CLIP models, fine-tuned on CIFAR10 using the respective source of randomness: random linear heads, data shuffling, and independent data subsets. Clearly, the use of random linear heads is the only method producing AGL behaviour (i.e. matching bias and slope of the two lines)

88 metrics [4, 23, 7, 20]. However, these approaches only provide upper bounds on the OOD error, and
 89 the bounds tend to be loose when evaluated on deep networks [25]. Another line of work looks at
 90 leveraging the model’s softmax predictions to predict the OOD performance [15, 14, 12, 10, 13].
 91 While these approaches show empirical promise in some settings, they are not expected to work in
 92 general and often fail in the presence of large shifts [12].

93 **ACL and AGL** Baek et al. [2] propose AGL, a recent approach for estimating OOD performance,
 94 that outperforms prior approaches across a variety of shifts. It is based on an earlier intriguing
 95 observation from [25, 30, 31, 32, 40, 35, 24]—there is a strong linear correlation between the
 96 probit-scaled ID and OOD performances of models across many distribution shift benchmarks (ACL).

97 Interestingly, Baek et al. [2] observes that when where ACL holds, the probit-scaled *agreement*
 98 *between models* is also strongly correlated and observe the *same* slope and bias. Furthermore, when
 99 accuracies do not show a linear correlation, agreements also do not. This phenomenon was called
 100 “agreement-on-the-line” (AGL).

101 Formally, given a pair of models f_1 and f_2 that map inputs to labels, accuracy and agreement can be
 102 defined as

$$\text{Acc}(f_1) = \mathbb{E}_{x,y \sim \mathcal{D}}[\ell(f_1(x), y)], \quad \text{Agr}(f_1, f_2) = \mathbb{E}_{x,y \sim \mathcal{D}}[\ell(f_1(x), f_2(x))], \quad (1)$$

103 where ℓ is the appropriate performance metric of interest. Note that while accuracy requires access to
 104 the labels y , agreement only requires access to unlabeled data and a pair of models. Thus, one can
 105 compute this line using OOD unlabeled data, and then estimate the OOD performance by linearly
 106 transforming the ID performance measured on ID validation data. See Appendix 5.3 for formal ALine
 107 methods to use AGL for OOD estimation.

108 **Training from scratch vs fine-tuning** A crucial component for AGL is the *diversity* of the ensemble
 109 predictions over which agreements are evaluated. If the models are not diverse enough, AGL is bound
 110 to fail. As an extreme, consider an ensemble of effectively identical models. Their ID and OOD
 111 agreement will always be 1, and there is no linear fit to estimate. Prior work on AGL has exclusively
 112 focused on training from scratch for several epochs, a very different regime from light fine-tuning. In
 113 this work, we focus on how to introduce sufficient diversity during *just* the fine-tuning process which
 114 can start from the *same* base FM and usually involves far fewer gradient steps.

115 3 Experiments and Results

116 **Fine-tuning.** In this work, we consider **linear probing (LP)** and **full fine-tuning (FFT)**. For LP,
 117 given features B_θ from the base model B , we train just the linear head v on top of frozen features
 118 such that the final classifier maps the score $v^\top B_\theta(x)$ to a predicted class. We refer to v as either a
 119 linear probe (classification) or span prediction head (QA). For FFT, we attach a linear head v and

Table 1: The MAPE (%) of predicting OOD performance using ALine and other baseline methods. Evaluations on QA tasks (SQuAD-Shifts) are performed over a set of models finetuned from multiple base FMs (LLaMa, GPT, OPT). Evaluations on the image classification datasets are conducted with CLIP models fine-tuned with linear probing.

OOD Dataset	ALine-D	ALine-S	Naive Agr	ATC	AC	DF
SQuAD-Shifts (averaged across 4 shifts)	1.68	2.55	19.48	9.16	45.04	4.54
CIFAR10C (averaged across shifts)	6.99	6.92	44.33	31.28	48.66	32.79
CIFAR10.1 (averaged across v4, v6)	2.42	3.03	41.52	6.48	54.57	8.51
CIFAR100C (averaged across shifts)	11.94	12.67	46.13	18.69	80.81	37.36
ImageNet V2 (averaged across 3 format)	4.96	5.03	47.65	8.96	77.34	7.86
WILDS (averaged across 3 benchmarks)	11.52	12.91	50.12	21.73	42.18	27.54

120 optimize the suitable loss function, but we *update all parameters* of the backbone such that the
 121 feature extract B_ϕ is updated. When infeasible to update all parameters natively, we perform *low-rank*
 122 *adaptation* (LoRA) [16] which uses trainable rank decomposition matrices to reduce the number of
 123 trainable parameters while still effectively updating the feature extractor B_ϕ . In this work, we do not
 124 distinguish between LoRA and FFT as they conceptually achieve the same effect, and show similar
 125 empirical trends in our studies. Refer to Appendix 5.1 for details on fine-tuned models and Appendix
 126 5.2 for specific fine-tuning parameters.

127 **Datasets** We study AGL for the tasks of QA and Image Classification. For QA, we fine-tune
 128 on the SQuAD v1.1 dataset [29] and evaluate on four distribution shifts present in SQuAD-Shifts
 129 (New Wiki, New York Times, Amazon, and Reddit) [24]. For image classification, we fine-tune
 130 on CIFAR10 [19], and then test on CIFAR10C [14], a dataset with 19 corruptions, some natural
 131 (Snow), and some synthetic (JPEG compression). We also test on the CIFAR10.1 dataset [30], which
 132 contains newer images for the same labels. We repeat the same for CIFAR100 [18], ImageNet-1k [33].
 133 We additionally validate our finding by testing on three natural shifts from the WILDS benchmark
 134 (FMoW, iWildCam, Camelyon17) [17].

135 3.1 Predicting OOD performance: single base foundation model

136 Consider the case where we have a *single* base FM to fine-tune. An overriding concern when
 137 calculating agreement is that even some randomness in the fine-tuning process may not be enough to
 138 overcome the underlying similarities in predictions due to the same base FM. To address this problem,
 139 we evaluate three possible methods for introducing diversity in the fine-tuning, to see what approach
 140 (if any) can lead to AGL behavior:

- 141 1. **Random linear heads.** Before fine-tuning, we initialize the last layer of the network (i.e.,
 142 the linear head) randomly, instead of via some zero-shot or pre-specified manner.
- 143 2. **Data shuffling.** We present the same data to each model, but shuffle the order for the data
 144 differently within each fine-tuning optimization run.
- 145 3. **Data subsetting.** We fine-tune each model with an independently sampled subset of the ID
 146 data. All models are trained on subsets of the same size.

147 Note that we perturb only one source of diversity at a time. For example, in the random linear head
 148 setting, all models start with a different initialization, but the data used for training is the same and
 149 seen in the same order. In the data shuffling setting, all models start with an identical arbitrary
 150 initialization, but the data used for training is seen in different orders; and so on.

151 For our study of image classification, we train a linear probe atop of CLIP, specifically the ViT-B/32
 152 model trained on LAION-2B [34]. For QA, we evaluate a collection of 50 fine-tuned models, all
 153 obtained by fine-tuning from the same checkpoint of a GPT2-Medium. We repeat the same procedure
 154 for OPT and BERT architectures, the details of which can be found in the Appendix (Sections 5.1
 155 and 5.5).

156 For the case of training models from scratch, it is well established that independent data subsetting
157 tends to lead to the greatest diversity of classifiers [26]. Nonetheless, in this setting we find rather
158 surprisingly, that model pairs trained with different *randomly initialized linear heads* achieve the
159 lowest OOD agreement for the same ID agreement. In fact, the ID versus OOD agreement matching
160 the slope of ID versus OOD accuracy. On the other hand, data ordering and data shuffling observe ID
161 versus OOD agreement that lies closer to the diagonal $y = x$ and away from the accuracy linear fit.
162 We show that this finding persists over numerous models and tasks.

163 3.2 Predicting OOD performance: multiple base foundation models

164 When multiple base foundation models (pretrained on different data) are accessible, it is unclear if
165 models with different bases would lie on different or similar accuracy lines, even if fine-tuned on the
166 same ID data. We observe that for certain extractive QA shifts, foundation models fine-tuned from a
167 wide range of base models *exhibit both ACL and AGL* (See Appendix 5.6 for details)

168 3.3 Results

169 Figure 1 shows the ACL/AGL trends for linear probes trained on top of CIFAR10 CLIP representa-
170 tions. One may suspect that such linear models would agree highly and AGL may break. However,
171 we see that contrary to the findings of Baek et al. [2], even linear models, when on top of neural
172 network features with the *right type of diversity*, may exhibit AGL. Interestingly enough, for the other
173 sources of diversity, we observe ACL and strongly linearly correlated agreement, but the latter at a
174 much higher rate OOD. We refer the reader to Appendix 5.9 for a more exhaustive evaluation. The
175 same observations, however not as stark, can be made for the fine-tuned LLMs. We refer the reader
176 to Appendix 5.5 to observe these trends on all four shifts within the SQuAD-Shifts dataset.

177 When considering multiple base foundation models, we first observe that base LLMs pretrained on
178 different corpora also lead to fine-tuned models that exhibit ACL. This is in contrast to the findings of
179 previous works [28, 35]. Second, the ID versus OOD agreement for pairs of models in this ensemble,
180 including pairs of different base foundation models, retains a strong linear correlation and the slope
181 and bias closely matches that of accuracy. As a result, different pretraining does not break AGL.

182 Table 1 shows the averaged MAPE (Mean Absolute Percentage Error) as calculated using the ALine
183 algorithm and other baseline methods for some dataset shifts (the full version for all datasets can be
184 found in Appendix 5.8). The QA ensembles are generated by fine-tuning multiple foundation models,
185 and the image classification ones are all CLIP linear-probes. Since AGL is demonstrated to hold well
186 in all these ensembles, the ALine MAE is able to surpass other methods; thus lending support to our
187 method to get AGL to hold for lightly fine-tuned models, and using it to estimate OOD performance.

188 4 Conclusion

189 We develop methods for extending AGL to lightly fine-tuned FMs to enable OOD performance
190 prediction in this emerging paradigm. We found that applying AGL directly may sometimes fail,
191 and proper utilization of this phenomena requires a careful tuning of the distribution of models in an
192 ensemble for their errors to be uncorrelated. Unlike the original paradigm of AGL, where models
193 observed tens or hundreds of epochs of training on the in-distribution dataset, we find that stochasticity
194 in specific optimization choices, specifically random initialization, is crucial for observing AGL in
195 lightly fine-tuned FMs. Second, though Baek et al. [2] posed AGL as a model centric phenomena
196 that is specifically only observed in neural network ensembles, we find that linear models can also
197 observe AGL when the data and the distribution shift contain certain structures (as is possible in the
198 CLIP representation space).

199 Our conclusion on AGL also sheds light on ACL (i.e. accuracy-on-the-line) in the presence of
200 foundation models, a phenomenon that is of independent interest. Some recent works have studied
201 the effect of different forms of fine-tuning on ACL [28, 1]. The main finding reported is that different
202 forms of fine-tuning lead to different slopes in the linear correlations, a term that is often called
203 “effective robustness”. In our results, we find that when fine-tuned the same way, models obtained
204 from *different base foundation models* all lie on the *same* line. This is particularly intriguing because
205 it goes against the common wisdom that the amount of pretraining data determines the effective
206 robustness. We leave these questions for future analysis.

207 **References**

- 208 [1] Anas Awadalla, Mitchell Wortsman, Gabriel Ilharco, Sewon Min, Ian Magnusson, Hannaneh
209 Hajishirzi, and Ludwig Schmidt. Exploring the landscape of distributional robustness for
210 question answering models. *arXiv preprint arXiv:2210.12517*, 2022.
- 211 [2] Christina Baek, Yiding Jiang, Aditi Raghunathan, and J Zico Kolter. Agreement-on-the-line:
212 Predicting the performance of neural networks under distribution shift. *Advances in Neural
213 Information Processing Systems*, 35:19274–19289, 2022.
- 214 [3] Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn.
215 The pushshift reddit dataset. In *Proceedings of the international AAAI conference on web and
216 social media*, volume 14, pages 830–839, 2020.
- 217 [4] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations
218 for domain adaptation. *Advances in neural information processing systems*, 19, 2006.
- 219 [5] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von
220 Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the
221 opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- 222 [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
223 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
224 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 225 [7] Corinna Cortes, Yishay Mansour, and Mehryar Mohri. Learning bounds for importance
226 weighting. *Advances in neural information processing systems*, 23, 2010.
- 227 [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
228 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
229 2018.
- 230 [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
231 deep bidirectional transformers for language understanding, 2019.
- 232 [10] Hady Elsahar and Matthias Gallé. To annotate or not? predicting performance drop under
233 domain shift. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Lan-
234 guage Processing and the 9th International Joint Conference on Natural Language Processing
235 (EMNLP-IJCNLP)*, pages 2163–2173, 2019.
- 236 [11] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason
237 Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse
238 text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 239 [12] Saurabh Garg, Sivaraman Balakrishnan, Zachary C Lipton, Behnam Neyshabur, and Hanie
240 Sedghi. Leveraging unlabeled data to predict out-of-distribution performance. *International
241 Conference on Learning Representations*, 2022.
- 242 [13] Devin Guillory, Vaishaal Shankar, Sayna Ebrahimi, Trevor Darrell, and Ludwig Schmidt. Pre-
243 dicting with confidence on unseen distributions. In *Proceedings of the IEEE/CVF international
244 conference on computer vision*, pages 1134–1144, 2021.
- 245 [14] Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common
246 corruptions and perturbations. In *7th International Conference on Learning Representations,
247 ICLR, 2019*.
- 248 [15] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
249 examples in neural networks. In *5th International Conference on Learning Representations,
250 ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- 251 [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
252 Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv
253 preprint arXiv:2106.09685*, 2021.

- 254 [17] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay
255 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al.
256 Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine*
257 *Learning*, pages 5637–5664. PMLR, 2021.
- 258 [18] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced
259 research). URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- 260 [19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
261 2009.
- 262 [20] Ilja Kuzborskij and Francesco Orabona. Stability and hypothesis transfer learning. In *Internat-*
263 *ional Conference on Machine Learning*, pages 942–950. PMLR, 2013.
- 264 [21] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-
265 image pre-training with frozen image encoders and large language models. *arXiv preprint*
266 *arXiv:2301.12597*, 2023.
- 267 [22] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
268 *arXiv:1711.05101*, 2017.
- 269 [23] Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation: Learning
270 bounds and algorithms. *arXiv preprint arXiv:0902.3430*, 2009.
- 271 [24] John Miller, Karl Krauth, Benjamin Recht, and Ludwig Schmidt. The effect of natural distri-
272 bution shift on question answering models. In *International conference on machine learning*,
273 pages 6905–6916. PMLR, 2020.
- 274 [25] John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal
275 Shankar, Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong
276 correlation between out-of-distribution and in-distribution generalization. In *International*
277 *Conference on Machine Learning*, pages 7721–7735. PMLR, 2021.
- 278 [26] Preetum Nakkiran and Yamini Bansal. Distributional generalization: A new kind of generaliza-
279 tion. *arXiv preprint arXiv:2009.08092*, 2020.
- 280 [27] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
281 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 282 [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
283 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
284 models from natural language supervision. In *International conference on machine learning*,
285 pages 8748–8763. PMLR, 2021.
- 286 [29] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions
287 for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- 288 [30] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do cifar-10
289 classifiers generalize to cifar-10? *arXiv preprint arXiv:1806.00451*, 2018.
- 290 [31] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet
291 classifiers generalize to imagenet? In *International conference on machine learning*, pages
292 5389–5400. PMLR, 2019.
- 293 [32] Rebecca Roelofs, Vaishaal Shankar, Benjamin Recht, Sara Fridovich-Keil, Moritz Hardt, John
294 Miller, and Ludwig Schmidt. A meta-analysis of overfitting in machine learning. *Advances in*
295 *Neural Information Processing Systems*, 32, 2019.
- 296 [33] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
297 Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-
298 Fei. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014. URL
299 <http://arxiv.org/abs/1409.0575>.

- 300 [34] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman,
301 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-
302 5b: An open large-scale dataset for training next generation image-text models. *Advances in*
303 *Neural Information Processing Systems*, 35:25278–25294, 2022.
- 304 [35] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig
305 Schmidt. Measuring robustness to natural distribution shifts in image classification. *Advances*
306 *in Neural Information Processing Systems*, 33:18583–18599, 2020.
- 307 [36] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
308 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas
309 Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes,
310 Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony
311 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian
312 Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut
313 Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,
314 Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta,
315 Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiao-
316 qing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng
317 Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien
318 Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation
319 and fine-tuned chat models, 2023.
- 320 [37] Trieu H Trinh and Quoc V Le. A simple method for commonsense reasoning. *arXiv preprint*
321 *arXiv:1806.02847*, 2018.
- 322 [38] Dequan Wang, Xiaosong Wang, Lilong Wang, Mengzhang Li, Qian Da, Xiaoqiang Liu, Xiangyu
323 Gao, Jun Shen, Junjun He, Tian Shen, et al. Medfmc: A real-world dataset and benchmark for
324 foundation model adaptation in medical image classification. *arXiv preprint arXiv:2306.09579*,
325 2023.
- 326 [39] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca
327 Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong,
328 et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF Conference on*
329 *Computer Vision and Pattern Recognition*, pages 7959–7971, 2022.
- 330 [40] Chhavi Yadav and Léon Bottou. Cold case: The lost mnist digits. *Advances in neural information*
331 *processing systems*, 32, 2019.
- 332 [41] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen,
333 Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam
334 Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke
335 Zettlemoyer. Opt: Open pre-trained transformer language models, 2022.
- 336 [42] Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba,
337 and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by
338 watching movies and reading books. In *Proceedings of the IEEE international conference on*
339 *computer vision*, pages 19–27, 2015.

340 5 Appendix

341 5.1 Models

342 **Extractive QA** We evaluate a collection of 125 fine-tuned models for our experiments in this
343 section. Each model is obtained by fine-tuning from the same checkpoint of a GPT2-Medium,
344 OPT-125M, and BERT. We individually present findings on both these families of models in the
345 following sections. Huggingface links to the base models we trained are in Appendix 5.7.

346 **Image Classification** We use CLIP [28], specifically the ViT-B/32 model trained on LAION-2B
347 [34] for our image classification tasks. Given its well-established 0-shot capabilities, a popular
348 method of fine-tuning CLIP for downstream tasks is to simply employ linear probing on top of the
349 CLIP representation. Thus, we are interested in evaluating the OOD performance of an ensemble of
350 models where the only difference is the linear head.

351 **Multiple Models** We train 41 models on the extractive QA benchmark SQuAD as in the previous
352 section, and observe their OOD performance to SQuAD-Shifts. We fine-tune OPT-125M [41], OPT-
353 350M, OPT-1.3B, GPT2-XL, GPT2-Large, GPT2-Medium, GPT2 [27], GPT-Neo-135M, Llama2-7B
354 [36], Alpaca-7B, and Vicuna-7B to extractive QA. OPT was pretrained on a wide variety of data
355 including BookCorpus [42], Stories [37], a subset of PILE [11], CCNews v2 corpus, and PushShift.io
356 Reddit [3]. Similarly, GPT2 was pretrained on BookCorpus while GPT-Neo was trained on PILE.
357 Llama2 was trained on an undisclosed set of publicly available data. Sprouting from Llama2,
358 Alpaca is additionally trained from Llama2 on instruction-following demonstrations while Vicuna is
359 additionally trained from Llama2 on user-shared conversations from ShareGPT.

360 5.2 Finetuning Specifics

361 We state here the specific parameters used in finetuning GPT2-Medium for extractive QA and CLIP
362 for image classification. Across the four different sources of diversity, the epochs are varied regardless
363 of the experiment. We train with AdamW as the optimizer [22]. For randomly initializing linear
364 heads we vary the seed for the head and keep all other values fixed. For changing the finetuning
365 hyperparameters, we vary the learning rate and weight decay. To shuffle the data, we change the data
366 seed that control the data ordering during training. And finally for data subsetting, we get different
367 proportions of the dataset which are independently sampled.

368 For the GPT2-Medium models, we train a total of 50 models for studying the sources of diversity.
369 For the CLIP models, we fine-tune upwards of 200 models (i.e. linear heads on top of the CLIP
370 representation) for the different vision datasets.

Table 2: Finetuning specifics for extractive QA (LR: learning rate, WD: weight decay, LS: linear head initialization seed, DS: data shuffling seed, DP: data subsetting proportion, EP: epochs, B: batch size)

Source of Diversity	GPT2-Medium	
	Varied	Fixed
Random linear heads	LS: varied	LR: 3×10^{-6} WD: 2×10^{-4} DS: fixed DP: 20% EP: 0-3 B: 4
Data shuffling	DS: varied	LR: 4×10^{-6} WD: 1×10^{-4} LS: fixed DP: 10% EP: 0-3 B: 4
Data subsetting	DP: 4.5% – 50%	LR: 2×10^{-6} WD: 1×10^{-4} DS: varied LS: fixed EP: 1 B: 4

Table 3: Finetuning specifics for OPT-125M (LR: learning rate, WD: weight decay, LS: linear head initialization seed, DS: data shuffling seed, SS: random subsetting seed, EP: epochs)

Source of Diversity	OPT-125M	
	Varied	Fixed
Random linear heads	LS: varied	LR: 4×10^{-7} WD: 1×10^{-5} DS: fixed SS: fixed EP: 10
Data shuffling	DS: varied	LR: 4×10^{-7} WD: 1×10^{-5} LS: fixed SS: fixed EP: 10
Data subsetting	SS: varied	LR: 4×10^{-7} WD: 1×10^{-5} DS: varied LS: fixed EP: 10

Table 4: Finetuning specifics for BERT (LR: learning rate, WD: weight decay, LS: linear head initialization seed, DS: data shuffling seed, SS: random subsetting seed, EP: epochs)

Source of Diversity	BERT	
	Varied	Fixed
Random linear heads	LS: varied	LR: 2×10^{-7} WD: 1×10^{-5} DS: fixed SS: fixed EP: 10
Data shuffling	DS: varied	LR: 2×10^{-7} WD: 1×10^{-5} LS: fixed SS: fixed EP: 10
Data subsetting	SS: varied	LR: 2×10^{-7} WD: 1×10^{-5} DS: varied LS: fixed EP: 10

Table 5: Finetuning specifics for CLIP (LR: learning rate, WD: weight decay, LS: linear head initialization seed, DS: data shuffling seed, DP: data subsetting proportion, EP: epochs, B: batch size)

Source of Diversity	CLIP + ViT-B/32 (LAION-2B)	
	Varied	Fixed
Random linear heads	LS: varied	LR: different per dataset WD: 0 DS: fixed DP: 100% EP: 1–100 B: 1024
Data shuffling	DS: varied	LR: different per dataset WD: 0 LS: fixed DP: 100% EP: 1–100 B: 1024
Data subsetting	DP: 10% – 50%	LR: different per dataset WD: 0 DS: varied LS: fixed EP: 1–100 B: 1024

371 **5.3 ALine-S/D**

372 ALine is the OOD accuracy estimating metric that utilizes AGL [2]. There are two methods within
 373 ALine: ALine-S and ALine-D

374 Given $Acc_{ID}(f_1)$ and $Agr_{OOD}(f_1, f_2)$, when agreement holds, the relationship between the agree-
 375 ment line and accuracy line is as follows.

$$\Phi^{-1}(Acc_{OOD}(f_1)) = a \cdot \Phi^{-1}(Acc_{ID}(f_1)) + b \Leftrightarrow \Phi^{-1}(Agr_{OOD}(f_1, f_2)) = a \cdot \Phi^{-1}(Agr_{ID}(f_1, f_2)) + b \quad (2)$$

376 To find $Acc_{OOD}(f_2)$, we can estimate the slope a and bias b as follows and

$$\hat{a}, \hat{b} = \arg \min_{a, b \in \mathbb{R}} \sum_{i \neq j} \left(\Phi^{-1}(\hat{Agr}_{OOD}(h_i, h_j)) - a \cdot \Phi^{-1}(\hat{Agr}_{ID}(h_i, h_j)) - b \right)^2 \quad (3)$$

377 With \hat{a} and \hat{b} , we can find $Acc_{OOD}(f_2)$ with the estimator for the ID accuracy $Acc_{ID}(f_1)$. This
 378 method is called Aline-S.

379 A similar method, ALine-D, uses pointwise accuracies and agreement of the model of interest instead
 380 of estimating the entire agreement line. If the models of interest are h and h' , then the following
 381 holds.

$$\frac{1}{2} (\Phi^{-1}(Acc_{OOD}(h)) + \Phi^{-1}(Acc_{OOD}(h'))) = \frac{a}{2} (\Phi^{-1}(Acc_{ID}(h)) + \Phi^{-1}(Acc_{ID}(h'))) + \frac{b}{2} \quad (4)$$

382 With the fact that $b = \Phi^{-1}(Agr_{OOD}(h, h')) - a \cdot \Phi^{-1}(Agr_{ID}(h, h'))$, we have

$$\begin{aligned} & \frac{1}{2} (\Phi^{-1}(Acc_{OOD}(h)) + \Phi^{-1}(Acc_{OOD}(h'))) \\ &= \Phi^{-1}(Agr_{OOD}(h, h')) + a \cdot \left(\frac{\Phi^{-1}(Acc_{ID}(h)) + \Phi^{-1}(Acc_{ID}(h'))}{2} - \Phi^{-1}(Agr_{ID}(h, h')) \right) \end{aligned} \quad (5)$$

383 With the two unknowns, $Acc_{OOD}(h)$ and $Acc_{OOD}(h')$, and one equations we cannot find the unknowns.
 384 However, with more overlapping pairs, we can get the same number equations as variables and find
 385 the OOD accuracy of a model of interest.

386 **5.4 Sources of Diversity (Image Classification)**

387 Figure 2 shows the three sources of diversity for the “Pixelate” and “JPEG-Compression” shifts in
 388 the CIFAR 10C OOD dataset. Table 6 shows the ALine-D MAE (%) for image classification on
 389 CIFAR10C (average across all 19 shifts).

Table 6: ALine-D MAE and MAPE for CLIP linear probing on CIFAR10 image classification. Note that the reported MAE and MAPE is averaged across all 19 CIFAR10C evaluated shifts.

Source of Diversity	CIFAR10C MAPE (%)	CIFAR10C MAE (%)
Random linear heads	15.88	5.74
Data shuffling	74.16	22.61
Data subsetting	25.94	7.39

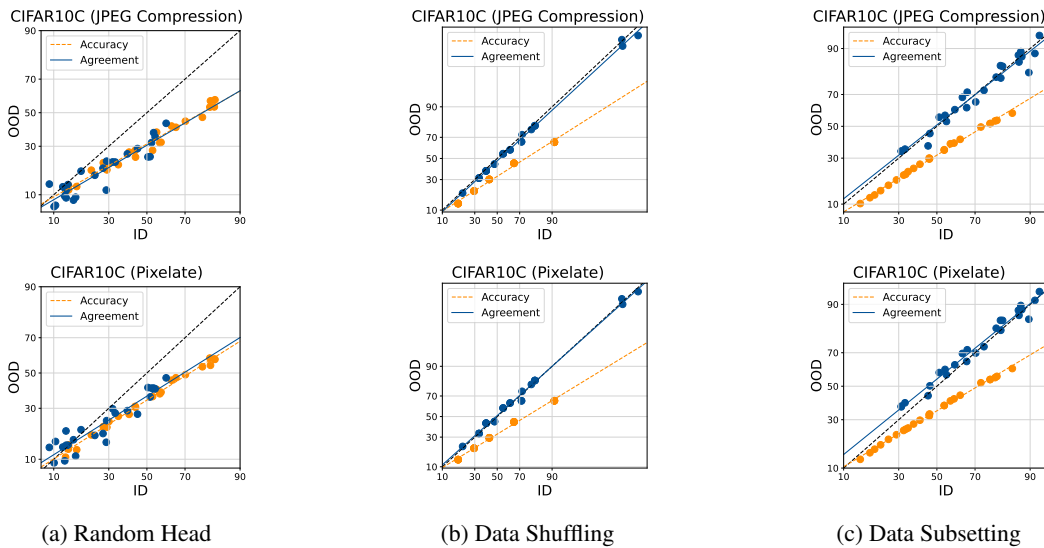


Figure 2: The ACL and AGL plots for the “JPEG Compression” (top row) and “Pixelate” (bottom row) fine-tuned using different sources of randomness

390 **5.5 Sources of Diversity (Question Answering)**

391 Figure 3 shows the three sources of diversity for all SQuAD-Shifts OOD datasets. Table 7 shows the
 392 ALine-D MAE for SQuAD-Shifts Amazon and Reddit.

Table 7: ALine-D MAPE(%) and MAE (%) on the SQuAD-Shifts Amazon and Reddit datasets when applied to sets of fully-finetuned models, trained using different sources of randomness

Source of Diversity	SQuAD-Shifts Amazon		SQuAD-Shifts Reddit	
	MAPE (%)	MAE (%)	MAPE (%)	MAE (%)
Random Linear Heads	6.34	0.69	3.48	0.79
Data Shuffling	10.30	4.18	9.59	4.32
Data Subsetting	16.21	5.2	13.94	4.71

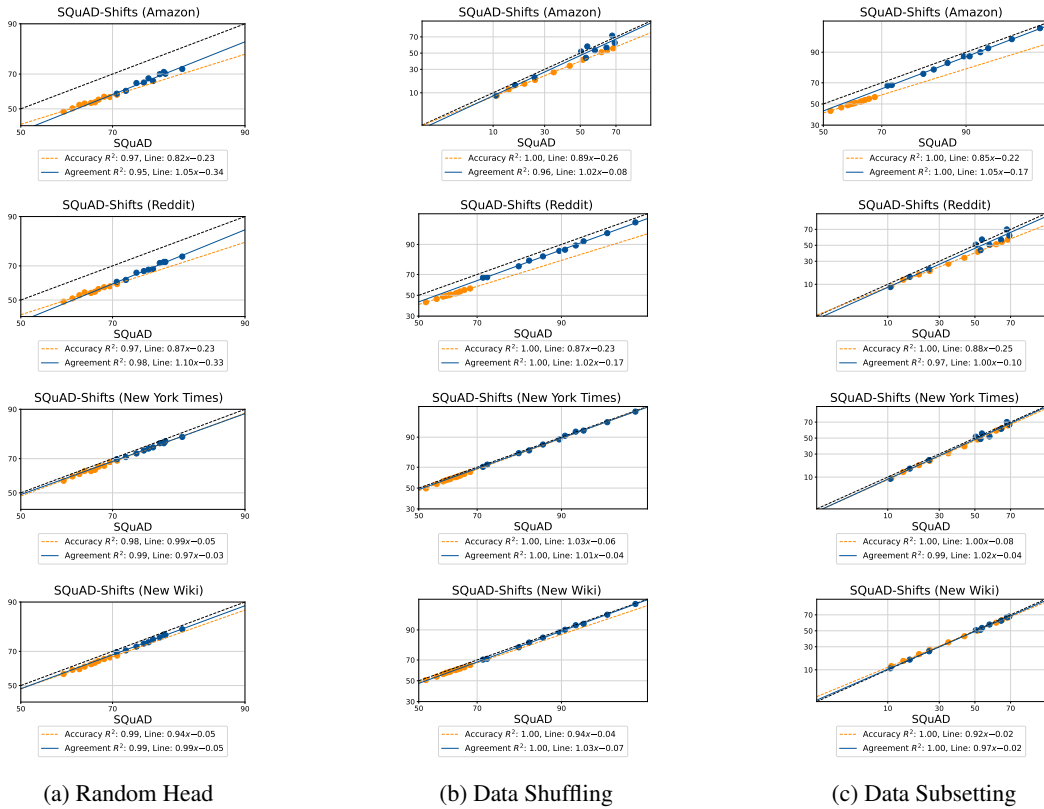


Figure 3: ID vs OOD trends of accuracy and agreement of LLMs finetuned for Question Answering from a single pretrained base model. Each column presents trends for different sources of stochasticity employed to obtain a diverse ensemble of finetuned models.

393 In this section, we also expand our evaluations to finetuned OPT-125M and BERT models for the
 394 extractive question answering task discussed in Section ???. For both of these base foundation models,
 395 we consider the three sources of diversity for finetuning i.e. using random linear heads, random
 396 ordering, and independent data subsetting, and plot the respective ID vs OOD accuracy of models
 397 and agreement between pairs of models in the resultant model set.

398 These experiments also afford us the chance to analyse the similarities and differences between the
 399 ACL/AGL trends exhibited by the model sets with GPT2-Medium, OPT-125M, and BERT as the
 400 base FM respectively. In particular, AGL is slightly worse for OPT-125M and BERT, and thus ALine
 401 has a higher error on OPT-125M and BERT than GPT2-Medium. However, we still see a consistent
 402 trend where AGL holds the best for random head initialization compared to data shuffling and data

403 subsetting; thus implying that the ALine error for random head initialization is the smallest out of all
 404 diversity sources. Thus, the importance of random head initialization applies to all models regardless
 405 of architecture in AGL.

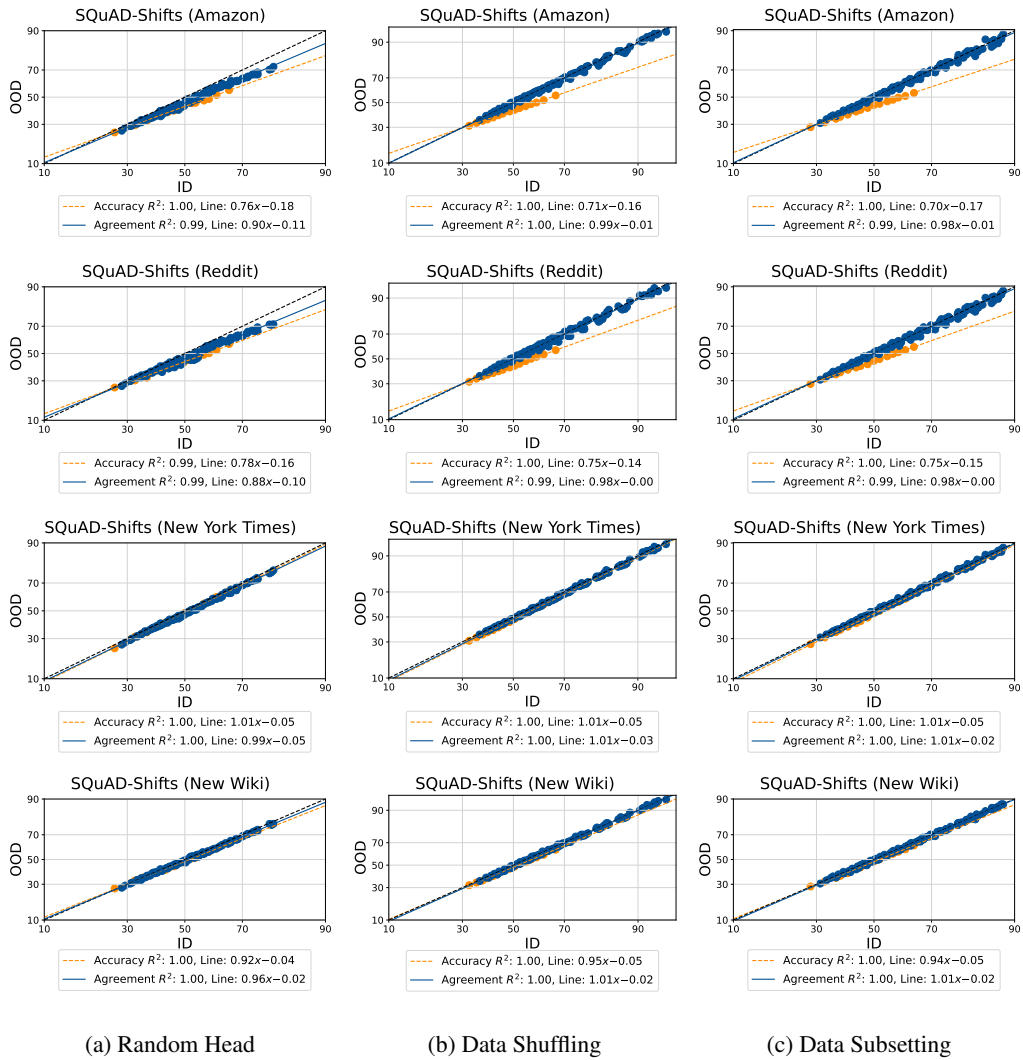


Figure 4: ID vs OOD trends of accuracy and agreement of LLMs finetuned for Question Answering from a single pretrained base model (OPT-125M). Similar to the GPT2-Medium results, these show that random linear head initialization is the best method to obtain model sets exhibiting AGL

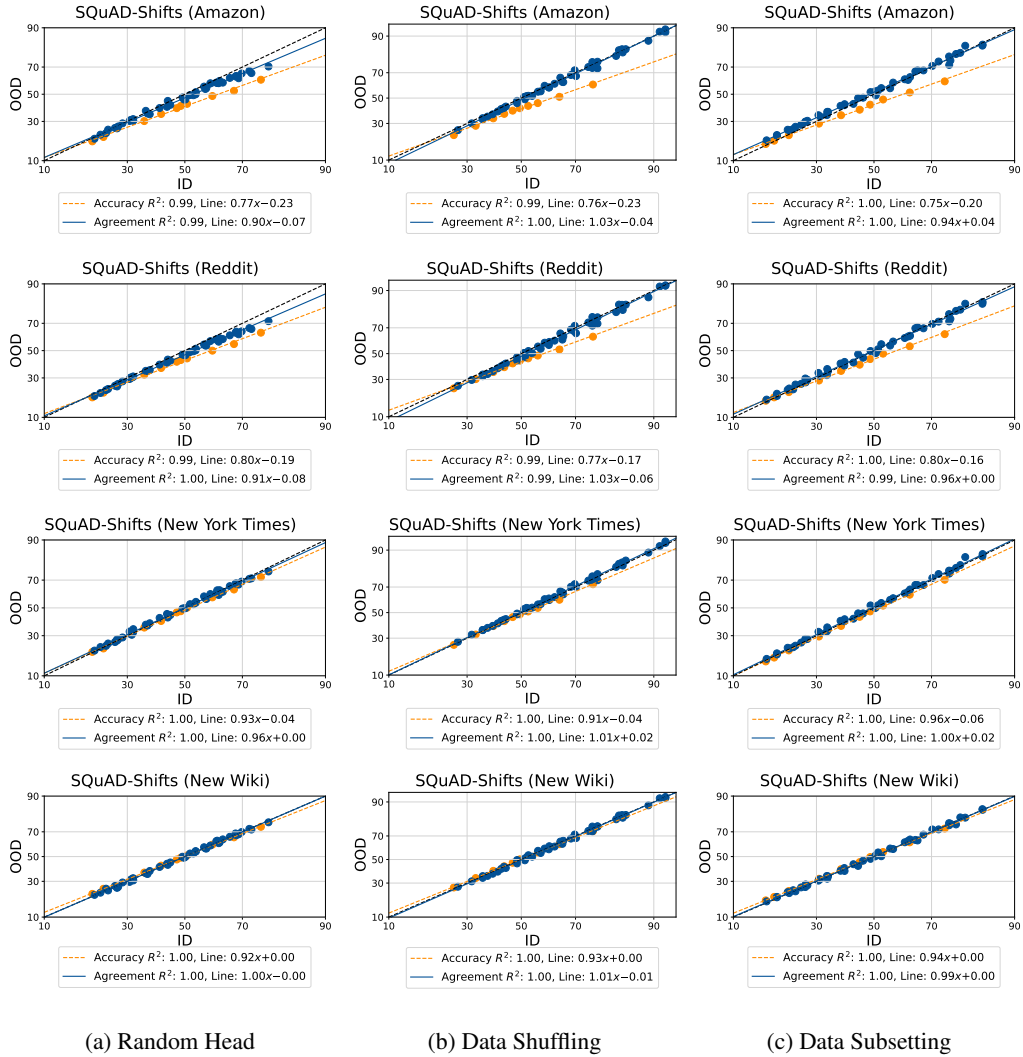


Figure 5: ID vs OOD trends of accuracy and agreement of LLMs finetuned for Question Answering from a single pretrained base model (BERT). Similar to the GPT2-Medium results, these show that random linear head initialization is the best method to obtain model sets exhibiting AGL

406 **5.6 Multiple Foundation Models**

407 Figure 6 shows AGL and ACL for different base models for all SQuAD-Shifts OOD datasets. We
408 have fine-tuned OPT-125M, OPT-350M, OPT-1.3B, GPT2-XL, GPT2-Large, GPT2-Medium, GPT2,
409 GPT-Neo-135M, Llama2-7B, Alpaca-7B, and Vicuna-7B. The links to the models are in Appendix
410 5.7.

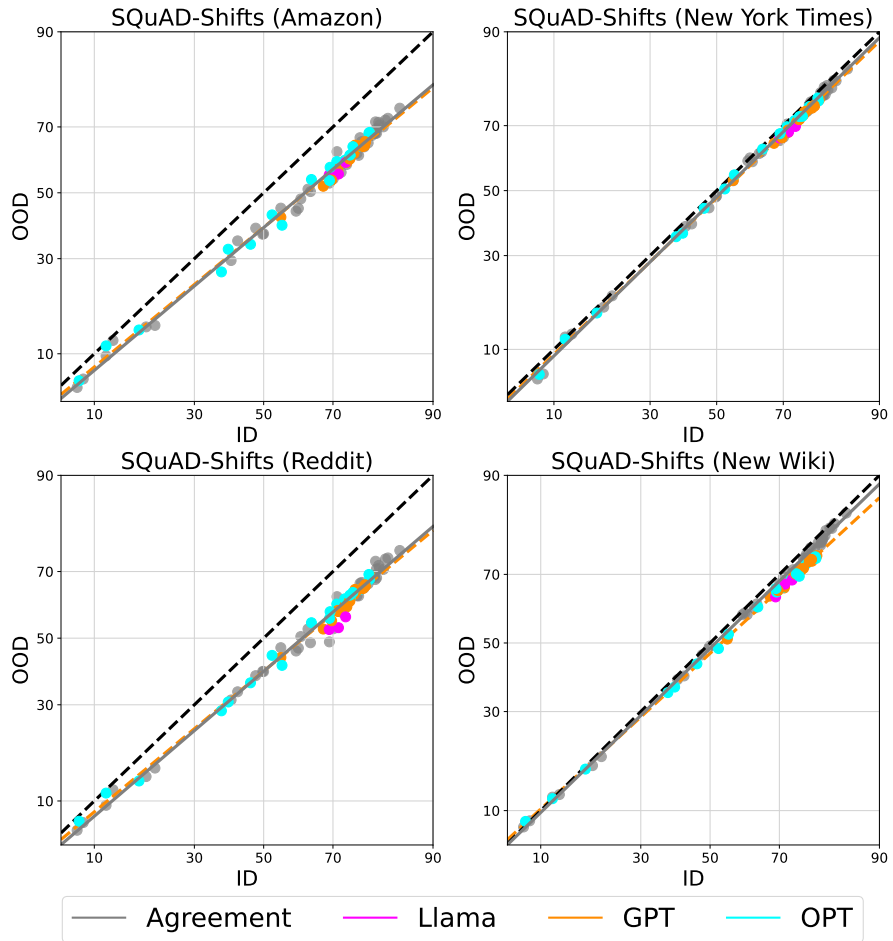


Figure 6: AGL when using different base models for SQuAD-Shifts

411 **5.7 Huggingface Links**

412 Here are the Huggingface links to the pretrained base foundation models we finetuned: GPT2 (<https://huggingface.co/gpt2>), GPT2-Medium (<https://huggingface.co/gpt2-medium>), GPT2-
413 Large (<https://huggingface.co/gpt2-large>), GPT2-XL (<https://huggingface.co/gpt2-xl>),
414 GPT-Neo-125M (<https://huggingface.co/EleutherAI/gpt-neo-125m>),
415 GPT-Neo-1.3B (<https://huggingface.co/EleutherAI/gpt-neo-1.3B>), OPT-125M
416 (<https://huggingface.co/facebook/opt-125m>), OPT-1.3B (<https://huggingface.co/facebook/opt-1.3b>),
417 Llama2-7B (<https://huggingface.co/meta-llama/Llama-2-7b-hf>),
418 Alpaca-7B (<https://huggingface.co/WeOpenML/Alpaca-7B-v1>),
419 Vicuna-7B (<https://huggingface.co/lmsys/vicuna-7b-v1.3>), BERT (<https://huggingface.co/bert-base-uncased>)
420
421

Table 8: The MAPE (%) of predicting OOD performance using ALine and other baseline methods. Evaluations on QA tasks (SQuAD-Shifts) are performed over a set of models finetuned from multiple base FMs (LLaMa, GPT, OPT). Evaluations on the image classification datasets are conducted with CLIP models fine-tuned with linear probing.

OOD Dataset	ALine-D	ALine-S	Naive Agr	ATC	AC	DF
SQuAD-Shifts Reddit	1.20	2.60	20.21	12.74	49.25	6.09
SQuAD-Shifts Amazon	1.64	3.10	20.40	15.35	51.06	7.39
SQuAD-Shifts Nyt	0.82	1.33	18.46	3.11	38.61	3.18
SQuAD-Shifts New Wiki	3.08	3.18	18.87	5.46	41.26	1.50
Average	1.68	2.55	19.48	9.16	45.04	4.54
CIFAR10C (averaged across shifts)	6.99	6.92	44.33	31.28	48.66	32.79
CIFAR10.1 (averaged across v4, v6)	2.42	3.03	41.52	6.48	54.57	8.51
CIFAR100C (averaged across shifts)	11.94	12.67	46.13	18.69	80.81	37.36
ImageNetC (averaged across shifts)	10.91	11.04	56.76	27.25	79.00	37.86
ImageNet V2 (averaged across 3 format)	4.96	5.03	47.65	8.96	77.34	7.86
fMoW-WILDS (val OOD split)	2.59	2.74	83.94	9.03	44.59	5.86
iWildCam-WILDS (val OOD split)	22.05	25.29	46.42	37.25	57.31	69.58
Camelyon17-WILDS (val OOD split)*	9.93	10.71	19.99	18.92	24.64	7.18

423 **5.9 Using Random-Head initialized fine-tuned CLIP models for other datasets**

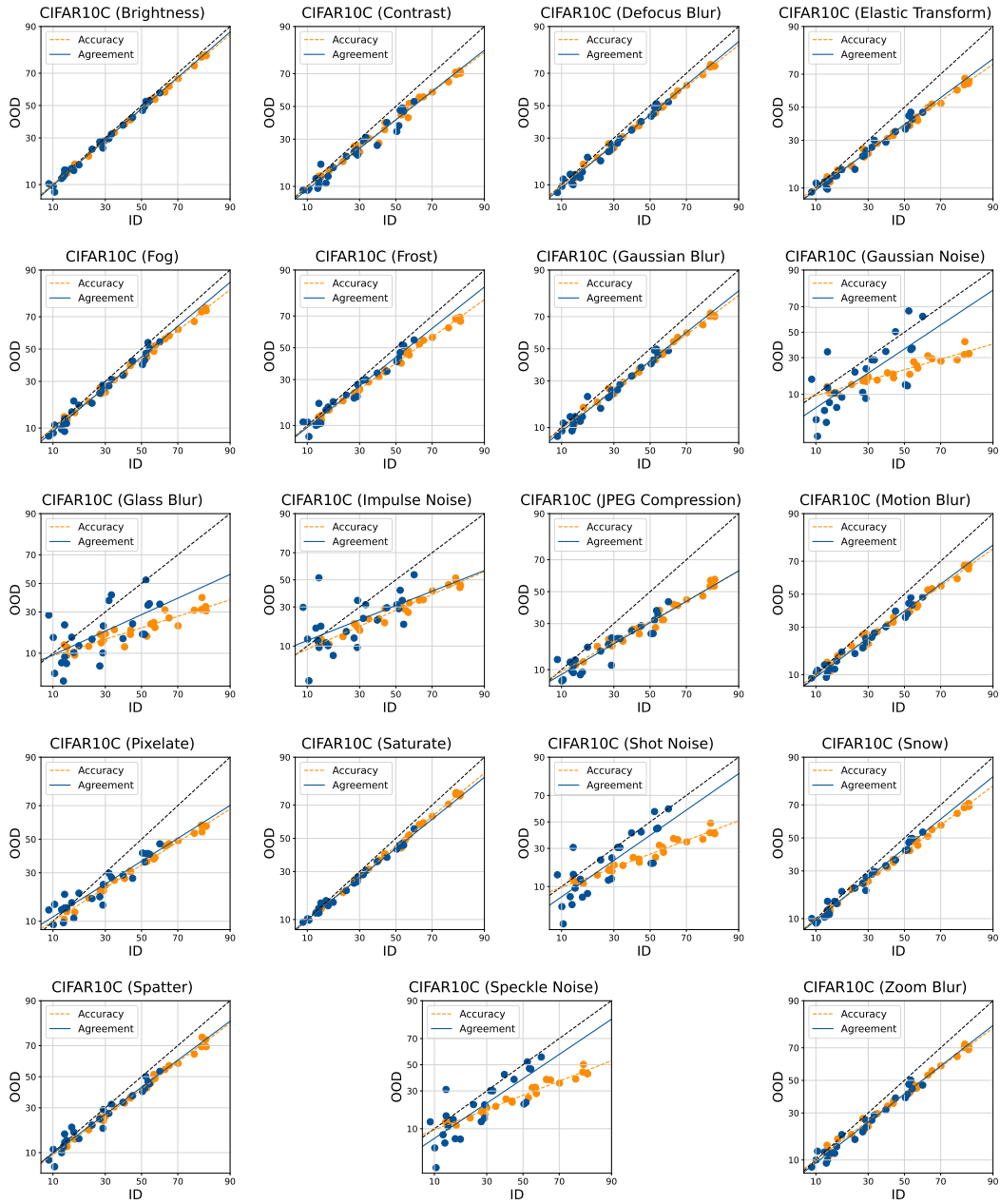


Figure 7: AGL and ACL for all C10C shifts with random head initialization fine-tuning.



Figure 8: AGL and ACL for the C10.1 shifts with random head initialization fine-tuning.

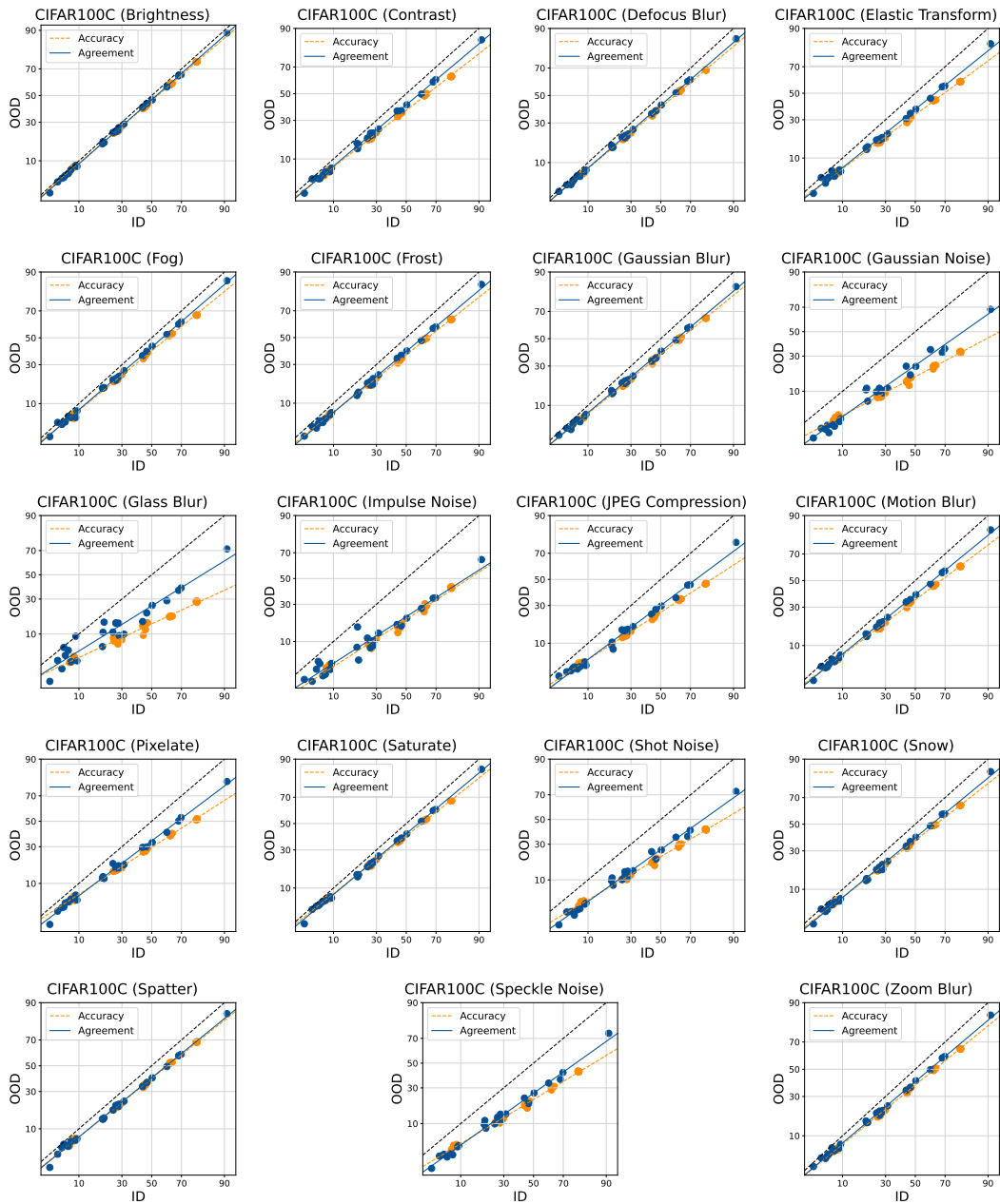


Figure 9: AGL and ACL for the C100C shifts with random head initialization fine-tuning.

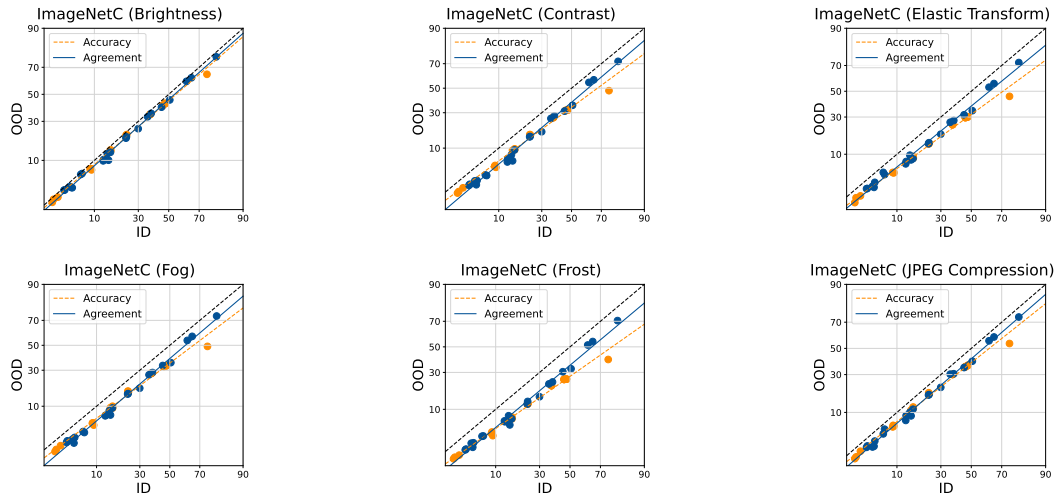


Figure 10: AGL and ACL for the ImageNetC shifts with random head initialization fine-tuning.

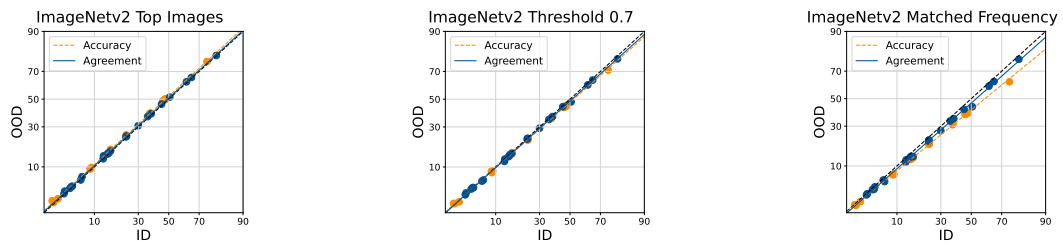


Figure 11: AGL and ACL for the ImageNet V2 shifts with random head initialization fine-tuning.

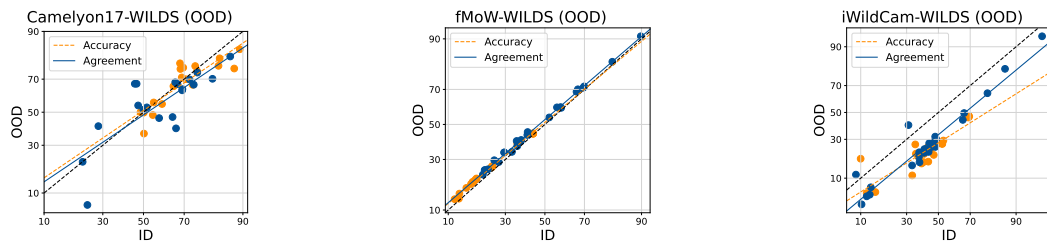


Figure 12: AGL and ACL for 3 benchmarks from the WILDS dataset with random head initialization fine-tuning.