

# On the Plasticity of Delta Parameters in Post-Trained Models

Anonymous ACL submission

## Abstract

Post-training has emerged as a crucial paradigm for adapting large-scale pre-trained models to various tasks, whose effects are fully reflected by delta parameters (i.e., the disparity between post-trained and pre-trained parameters). While numerous studies have explored delta parameter properties via operations like pruning, quantization, low-rank approximation, and extrapolation, a fundamental question remains: what properties of delta parameters are essential for maintaining performance? In this work, we investigate delta parameter properties along two dimensions: magnitude and sign. Through experiments on instruct language models, reasoning language models, and vision models, we find that delta parameters exhibit considerable *plasticity*: individual values, distribution shape, relative relationships, and even signs can be substantially modified while maintaining post-trained model’s performance. To understand these phenomena, we develop a loss-based theoretical framework that analyzes editing effects through a second-order Taylor expansion. Our analysis introduces the concept of editing intensity, which helps explain the stability boundaries of different editing operations, and identifies mean and relative relationships as key factors from a theoretical perspective.

## 1 Introduction

Post-training has become a critical step in developing large-scale models (Han et al., 2024; Xin et al., 2024; Dodge et al., 2020; Zhao et al., 2023). Through supervised fine-tuning and reinforcement learning, post-training endows pre-trained models with diverse capabilities such as instruction following (Rafailov et al., 2023; Ethayarajh et al., 2024), mathematical reasoning (Luo et al., 2023; Tong et al., 2024), code generation (Wang et al., 2025), and visual recognition (Chen et al., 2022; Sandler et al., 2022). The effect of post-training is fully reflected in the *delta parameters*, which

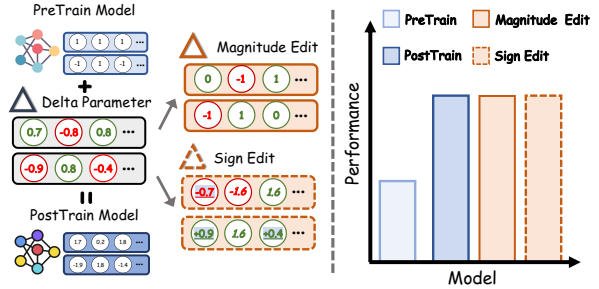


Figure 1: Delta parameters exhibit plasticity in both magnitude and sign. We investigate editing operations that modify magnitude (e.g., dropping and rescaling) or flip signs (with rescaling). Despite substantial modifications to delta parameters, the edited models can largely preserve the post-trained model’s performance.

are defined as the difference between post-trained and pre-trained parameters (Ilharco et al., 2023; Yu et al., 2024). Understanding the properties of delta parameters is therefore crucial for understanding post-training itself.

Recent years have witnessed various methods that edit delta parameters for different benefits. For instance, DARE (Yu et al., 2024) and DELLA-Merging (Deep et al., 2024) showed that models can achieve comparable performance with only a small fraction of delta parameters. BitDelta (Liu et al., 2024) demonstrated that delta parameters can be quantized to 1 bit with modest performance degradation. EXPO (Zheng et al., 2024) observed that extrapolating delta parameters with a suitable scaling factor can even enhance alignment performance. These works demonstrate that editing delta parameters can yield benefits ranging from efficient storage to improved alignment. However, they focus on different operations with different objectives, leading to scattered findings. A fundamental question remains unanswered: *What properties of delta parameters are essential for maintaining performance, and what can be freely manipulated?*

In this work, we systematically investigate delta parameter properties along two dimensions: *magni-*

069 *tude* (the absolute value) and *sign* (the direction of  
070 change). We conduct experiments across instruct  
071 language models (LLaMA-3-8B-Instruct (Dubey  
072 et al., 2024), Mistral-7B-Instruct-v0.3 (Jiang et al.,  
073 2023), Qwen2-7B-Instruct (Yang et al., 2024)),  
074 reasoning language models ( Qwen3-1.7B (Team,  
075 2025)), and vision models (ViT-B-32 (Radford  
076 et al., 2021)), covering post-training techniques  
077 including SFT, RLHF (Qwen et al., 2025), and  
078 RLVR (DeepSeek-AI et al., 2025). As shown in  
079 Figure 1, we find that delta parameters exhibit con-  
080 siderable *plasticity*: individual values, distribution  
081 shape, relative relationships, and even signs can  
082 be substantially modified while maintaining post-  
083 trained model’s performance. In the magnitude  
084 dimension, we find that within a reasonable edit-  
085 ing range, what matters more is the overall statisti-  
086 cal properties such as the mean of the magnitude,  
087 rather than individual parameter values. More sur-  
088 prisingly, in the sign dimension, we discover that a  
089 substantial proportion of signs can be flipped while  
090 still maintaining comparable performance to the  
091 post-trained model. This finding suggests that the  
092 direction of parameter updates, often assumed to  
093 be important in prior work (Yadav et al., 2023; Liu  
094 et al., 2024), also exhibits plasticity.

095 To understand these phenomena, we develop a  
096 loss-based theoretical framework. We analyze the  
097 effect of delta parameter editing through a second-  
098 order Taylor expansion of the loss function. This  
099 analysis reveals that the stability of editing opera-  
100 tions is related to an *editing intensity* term, which  
101 explains why high drop rates and sign-flip opera-  
102 tions are more prone to performance degrada-  
103 tion. Furthermore, our theoretical analysis iden-  
104 tifies mean and relative relationships as key factors  
105 affecting performance, providing a principled ex-  
106 planation for our empirical findings.

## 107 2 Preliminaries

### 108 2.1 Notation

109 Let  $W_{pre} \in \mathbb{R}^{d \times k}$  denote the parameters of a pre-  
110 trained model, where  $d$  and  $k$  represent the out-  
111 put and input dimensions. A post-trained model  
112 with parameters  $W_{post} \in \mathbb{R}^{d \times k}$  can be derived  
113 from the pre-trained backbone through supervised  
114 fine-tuning or reinforcement learning. The delta  
115 parameters are defined as the difference between  
116 post-trained and pre-trained parameters:  $\Delta W =$   
117  $W_{post} - W_{pre} \in \mathbb{R}^{d \times k}$ . Since delta parameters  
118 reflect the complete effect of post-training, under-

standing their properties is crucial for understand-  
ing post-training itself.

Delta parameter editing refers to applying a  
transformation  $\mathcal{F}$  to the original delta paramet-  
ers, yielding edited delta parameters  $\Delta \widetilde{W}_{edit} =$   
 $\mathcal{F}(\Delta W)$ . The final edited model is then obtained  
as  $W_{edit} = W_{pre} + \Delta \widetilde{W}_{edit}$ . Various editing oper-  
ations have been explored in prior work, including  
pruning, quantization, and extrapolation.

### 128 2.2 Representative Methods

129 DARE (Yu et al., 2024) is a representative delta  
130 parameter editing method designed to reduce pa-  
131 rameter redundancy and further mitigate conflicts  
132 in model merging. Specifically, DARE first drops  
133 delta parameters with probability  $p$ , then rescales  
134 the remaining parameters by  $1/(1-p)$ :

$$135 \Delta \widetilde{W}_{DARE} = \frac{1-M}{1-p} \odot \Delta W \quad (1)$$

136 where  $M \sim \text{Bernoulli}(p)$  is a random binary mask  
137 and  $\odot$  denotes element-wise multiplication. With  
138 this operation, DARE can drop up to 90% of delta  
139 parameters while maintaining model performance.

140 BitDelta (Liu et al., 2024) proposes a quantiza-  
141 tion method for delta parameters. It preserves only  
142 the sign  $\text{sign}(\Delta W)$  and replaces all magnitudes  
143 with the average magnitude  $\text{AVG}(|\Delta W|)$ :

$$144 \Delta \widetilde{W}_{BitDelta} = \text{AVG}(|\Delta W|) \cdot \text{sign}(\Delta W) \quad (2)$$

145 In this way, BitDelta quantizes delta parameters to  
146 1-bit while maintaining most of the model perfor-  
147 mance with slight degradation.

148 These two methods demonstrate that aggressive  
149 modifications to delta parameters do not necessarily  
150 cause severe performance degradation. This raises  
151 a natural question: what properties of delta pa-  
152 rameters are essential for maintaining post-trained  
153 model performance? In the following sections, we  
154 systematically investigate this question along two  
155 dimensions: magnitude and sign.

## 156 3 Plasticity of Delta Parameters

157 In this section, we systematically investigate the  
158 properties of delta parameters through experi-  
159 ments. We conduct experiments on instruct lan-  
160 guage models (LLaMA-3-8B-Instruct (Dubey et al.,  
161 2024), Mistral-7B-Instruct-v0.3 (Jiang et al., 2023),  
162 Qwen2-7B-Instruct (Yang et al., 2024)), reasoning  
163 language models ( Qwen3-1.7B (Team, 2025)), and  
164 vision models (ViT-B-32 (Radford et al., 2021)).

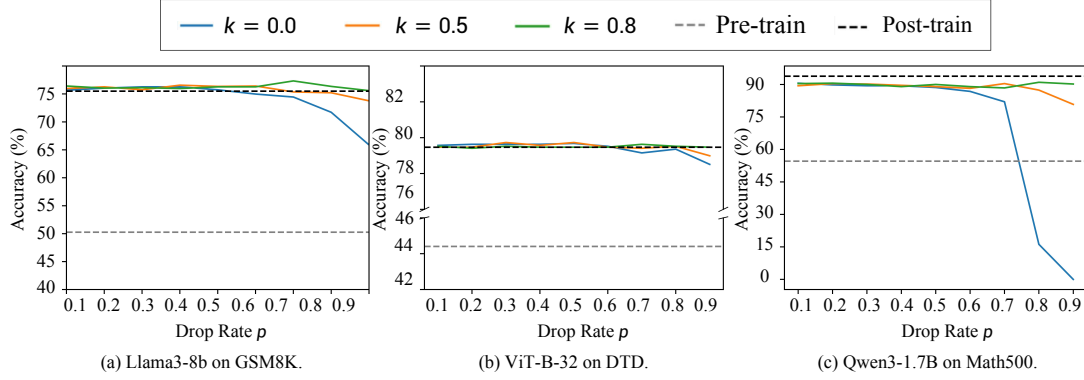


Figure 2: Performance under magnitude editing with varying drop rate  $p$  and scaling coefficient  $k$ .

These models cover most post-training techniques, including SFT, RFT, RLHF, and RLVR. We select appropriate evaluation tasks for each category of models. For instruct language models, we evaluate on 8 tasks: ARC Challenge (Clark et al., 2018), GSM8K (Cobbe et al., 2021), HellaSwag (Zellers et al., 2019), HumanEval (Chen et al., 2021), IFEval (Zhou et al., 2023), MMLU (Hendrycks et al., 2020), TruthfulQA (Lin et al., 2021), and Winogrande (Sakaguchi et al., 2021). For reasoning language models, we evaluate on MATH-500 (Lightman et al., 2023), AIME 2025 (American Invitational Mathematics Examination problems), GPQA Diamond (Rein et al., 2024), and LiveCodeBench (Jain et al., 2024). For vision models, we evaluate on 8 image classification tasks: Cars (Krause et al., 2013), DTD (Cimpoi et al., 2014), EuroSAT (Helber et al., 2019), GT-SRB (Stallkamp et al., 2011), MNIST (LeCun et al., 2010), RESISC45 (Cheng et al., 2017), SUN397 (Xiao et al., 2016), and SVHN (Netzer et al., 2011).

### 3.1 Plasticity in Magnitude

To investigate the plasticity of delta parameters in magnitude, we begin with DARE, a representative delta parameter editing method. DARE randomly selects a proportion  $p$  of delta parameters and sets them to zero, then rescales the remaining parameters by  $1/(1-p)$ . With this operation, DARE can drop up to 90% of delta parameters while maintaining model performance. The original paper explains this phenomenon through the lens of expected embeddings. Consider a linear transformation  $h = Wx + b$ , which is the basic operation in neural networks. Let  $\Delta W$  and  $\Delta b$  denote the delta parameters. After applying DARE with rescale factor  $\gamma$ , the expectation of the output becomes:

$$\mathbb{E}[\hat{h}] = W_{pre}x + b_{pre} + (1-p) \cdot \gamma \cdot (\Delta Wx + \Delta b)$$

By setting  $\gamma = 1/(1-p)$ , we have  $\mathbb{E}[\hat{h}] = h$ , i.e., the expected output is preserved. This preservation of expected output is argued to be the key to maintaining model performance.

DARE sets the selected parameters to zero (i.e., multiplies them by 0). A natural question arises: can we multiply by coefficients other than zero, scaling up or down some parameters while rescaling the rest, and still recover model performance? In other words, if we randomly select delta parameters with probability  $p$  and multiply them by a coefficient  $k$ , what should the rescale factor  $\gamma$  be for the remaining parameters? Based on DARE’s theoretical framework, the rescale factor for the remaining parameters should be  $(1-kp)/(1-p)$  to preserve the expected output (detailed derivation in Appendix A). This yields a generalized formulation:

$$\Delta \tilde{W} = k \cdot M \odot \Delta W + \frac{1-kp}{1-p} \cdot (1-M) \odot \Delta W \quad (3)$$

where  $M \sim \text{Bernoulli}(p)$  is a random binary mask. When  $k = 0$ , this reduces to the original DARE. When  $k = 1$ , no editing is performed.

To verify this hypothesis, we conduct experiments with different values of  $k$  and drop rates  $p$ . Figure 2 shows the results on representative models and datasets. When the drop rate  $p$  is relatively small, model performance remains nearly identical to the original post-trained model across a wide range of  $k$  values. When  $p$  is larger, performance slightly decreases but remains comparable to the original DARE setting ( $k = 0$ ). Similar patterns are observed on other settings (see Appendix B). These results indicate that scaling a subset of delta parameters by various coefficients, while appropriately rescaling the remaining parameters, can maintain nearly the same model performance.

The above findings demonstrate that delta parameters exhibit considerable plasticity in magnitude:

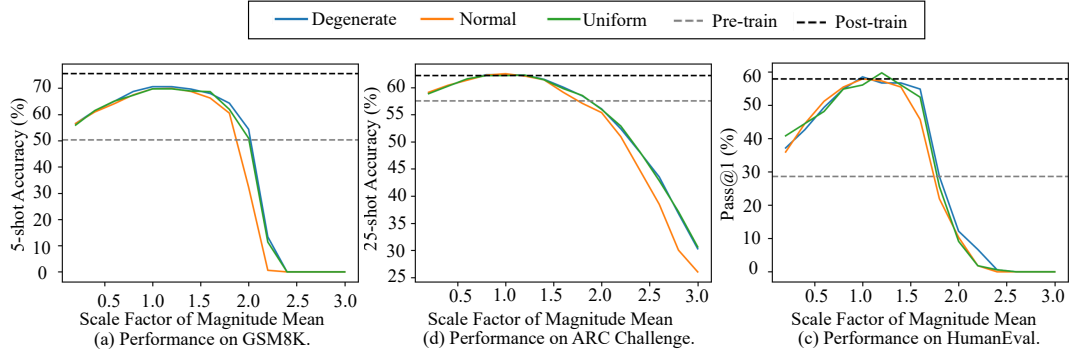


Figure 3: Performance under varying magnitude mean and distribution shape. The x-axis represents the scaling factor of the mean (1.0 is the original mean). Different curves correspond to different distribution shapes: uniform, normal, and degenerate.

we can scale a subset of parameters by different coefficients while rescaling the rest to preserve performance. This raises a further question: what properties of magnitude are truly essential for maintaining model performance? We consider three levels of properties, from fine-grained to coarse-grained: (1) the specific value of each individual parameter, (2) the relative relationships among parameters (e.g., ordering by magnitude), and (3) the global statistical properties (e.g., mean of the magnitude, distribution shape). We design a series of experiments to keep the sign and investigate the importance of each properties.

**Specific Values.** To investigate whether specific values are essential, we design an experiment that changes specific values while preserving relative relationships and global statistics. Specifically, we apply a power transformation to all magnitude values: raising each magnitude to the power of  $\alpha$  (we test  $\alpha = 0.5$  and  $\alpha = 1.5$ ), which alters each individual value but maintains the relative ordering among parameters. We then rescale the transformed magnitudes to restore the original mean. Table 1 shows the results. We find that under both transformations, performance remains nearly unchanged across all datasets. This suggests that specific magnitude values have limited impact on the capabilities learned through post-training.

**Relative Relationships.** DARE’s zero-out operation already suggests that relative relationships can be partially disrupted without severe performance degradation. To systematically investigate this factor, we design a shuffle experiment: we randomly shuffle a proportion  $r$  of delta parameter magnitudes across positions, varying  $r$  from 10% to 100%. This operation progressively destroys relative relationships while preserving the

Task	Original	Power 0.5	Power 1.5
ARC Challenge	62.20	62.46	61.95
GSM8K	75.51	74.67	75.36
HellaSwag	78.84	79.14	78.22
IFEval	47.12	47.28	47.07
MMLU	65.82	65.53	64.89
TruthfulQA	51.65	51.41	52.21
Winogrande	75.77	76.09	75.45

Table 1: Performance comparison between post-trained model and power & rescale model on LLaMA-3-8B-Instruct.

global distribution and the specific values. Figure 4 shows the results. When the shuffle rate is low, we observe limited performance degradation. As the shuffle rate increases, performance gradually decreases on some datasets such as GSM8K, but remains reasonable even at 100% shuffle rate. On other datasets such as ARC-Challenge, performance is almost unaffected even at high shuffle rates. This suggests that relative relationships contribute to performance to some extent, particularly when severely disrupted.

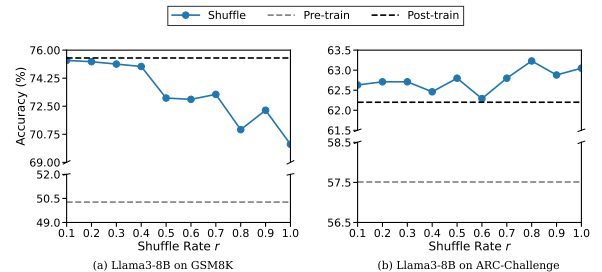


Figure 4: Performance of LLaMA-3-8B-Instruct on GSM8K and ARC Challenge with different shuffle rates.

**Global Statistical Properties.** We investigate two aspects of global statistics: the distribution shape and the mean of the magnitude ( $|\Delta W|$ ). We design an experiment that jointly varies both fac-

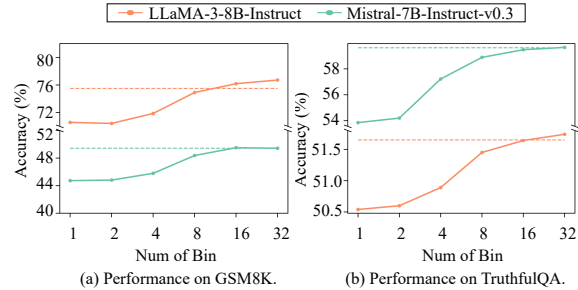
293 tors. For the distribution shape, we consider three  
 294 options: (1) a uniform distribution, (2) a normal  
 295 distribution, and (3) a degenerate distribution (i.e.,  
 296 all magnitudes set to the same value). For the mean  
 297 of the magnitude, we scale all magnitudes by a  
 298 constant factor  $\alpha$  ranging from 0.1 to 3.0. Figure 3  
 299 shows the results on three representative tasks. We  
 300 observe two clear patterns. First, at the same mean  
 301 value, the three distributions achieve nearly identi-  
 302 cal performance across all tasks. This indicates that  
 303 the specific distribution shape has limited impact  
 304 on performance. Second, when the mean deviates  
 305 from the original mean of magnitude, performance  
 306 degrades noticeably. These results indicate that the  
 307 mean of magnitude is a relatively important low-  
 308 dimensional indicator, while the distribution shape  
 309 has minimal impact.

310 The above experiments suggest a hierarchy of  
 311 importance among magnitude properties within our  
 312 experimental setting. The mean of the magnitude  
 313 appears to be the most sensitive factor. The re-  
 314 lative relationships among parameters have some  
 315 impact on performance, particularly when severely  
 316 disrupted. Within a reasonable editing range, the  
 317 specific value of each individual parameter and the  
 318 distribution shape show less sensitivity.

319 This understanding is consistent with the phe-  
 320 nomena observed in existing delta parameter edit-  
 321 ing methods. For DARE, although individual val-  
 322 ues are perturbed through random dropping and  
 323 rescaling, the mean of magnitude is preserved by  
 324 the rescale operation, and relative relationships are  
 325 partially maintained among the non-dropped pa-  
 326 rameters. This may explain why DARE can main-  
 327 tain performance. For BitDelta, all magnitudes are  
 328 replaced with the mean value, which preserves the  
 329 mean but completely destroys relative relationships.  
 330 According to our analysis, this would be expected  
 331 to cause some performance degradation, which is  
 332 consistent with the empirical observations in the  
 333 original paper.

334 Based on this understanding, we hypothesize  
 335 that partially restoring relative relationships could  
 336 improve BitDelta’s performance. To verify this, we  
 337 propose a simple modification: instead of replac-  
 338 ing all magnitudes with a single value, we parti-  
 339 tion parameters into  $K$  bins based on their original  
 340 magnitude ranking. Parameters within each bin  
 341 are then assigned the mean magnitude of that bin.  
 342 When  $K = 1$ , this reduces to the original BitDelta.  
 343 As  $K$  increases, more relative relationship infor-  
 344 mation is preserved. Figure 5 shows the results.

345 Performance improves consistently as  $K$  increases.  
 346 When  $K = 16$ , performance approaches that of the  
 347 original post-trained model. This supports our anal-  
 348 ysis: while the mean of the magnitude is important,  
 349 partially preserving relative relationships provides  
 additional benefits.



350 Figure 5: Effectiveness of increasing the number of bins  
 351 in BitDelta. The left subplot shows the performance  
 352 of LLaMA3-8B-Instruct and Mistral-7B-Instruct-v0.3  
 353 on the GSM8K dataset. The right subplot shows the  
 354 performance on the TruthfulQA dataset. In each subplot,  
 355 we use the dashed line to represent the performance of  
 356 the original post-trained model.

### 351 3.2 Plasticity in Sign

352 The previous subsection demonstrates that delta  
 353 parameters exhibit considerable plasticity in mag-  
 354 nitude. In this section, we investigate whether delta  
 355 parameters also exhibit plasticity in the sign dimen-  
 356 sion. Intuitively, the sign represents the direction  
 357 of parameter adjustment during post-training, in-  
 358 dicated whether a parameter should increase or  
 359 decrease relative to the pre-trained value. This  
 360 directional information is often assumed to be im-  
 361 portant (Yadav et al., 2023; Liu et al., 2024).

362 To investigate whether signs can be modified, we  
 363 extend the generalized formulation in Equation 3  
 364 to negative  $k$  values. When  $k < 0$ , the selected  
 365 parameters are multiplied by a negative coefficient,  
 366 which flips their signs and scales their magnitudes.  
 367 Then we rescale the remaining parameters by  $(1 -$   
 368  $kp)/(1 - p)$ , which is expected to preserve the  
 369 expected output.

370 To verify this, we conduct sign-flip experiments  
 371 across multiple models. Specifically, we con-  
 372 sider two representative settings:  $k = -0.5$  (flip-  
 373 ping signs while reducing magnitude by half) and  
 374  $k = -1.0$  (fully flipping signs without magnitude  
 375 change). For each setting, we vary the flip propor-  
 376 tion  $p$  and observe the resulting performance. Fig-  
 377 ure 6 shows representative the results(Full results  
 378 are shown in Appendix B). We observe interesting  
 379 patterns across different ranges of  $p$ . When the

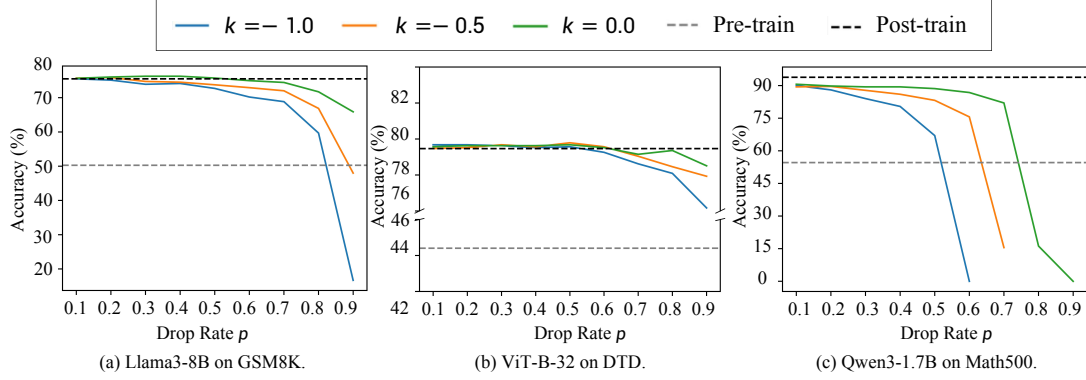


Figure 6: Performance under sign editing with varying flip rate  $p$  and scaling coefficient  $k < 0$ .

flip proportion is small, almost all models across all tasks can tolerate sign flipping with only minor performance degradation. This is a notable finding given the common assumption that signs encode essential directional information. When  $p$  becomes larger, different models and tasks exhibit varying degrees of robustness. For instance, LLaMA-3-8B-Instruct on GSM8K can tolerate up to 60% sign flipping while maintaining reasonable performance. More strikingly, ViT-B-32 on several vision tasks can tolerate up to 90% complete sign flipping with almost no performance degradation after rescaling. These results suggest that the plasticity of signs varies across models and tasks, but a substantial degree of sign modification is generally tolerable.

Combined with our findings on magnitude, delta parameters show plasticity in both magnitude and sign: within a reasonable editing range, individual values, distribution shape, relative relationships, and even signs can be substantially modified while maintaining model performance.

## 4 Understanding the Plasticity of Delta Parameters

In the previous section, we observed that delta parameters exhibit considerable plasticity in both magnitude and sign. At the same time, we also observed some stability boundaries: performance degrades sharply when the drop rate is too high, and sign-flip becomes unstable earlier than magnitude-only editing at comparable modification rates. DARE’s theoretical framework provides a useful intuition by approximately preserving expected outputs, but it does not readily explain these phenomena. In this section, we enrich this understanding by analyzing the loss change induced by editing perturbations using a second-order surrogate, which allows us to better understand the plasticity of delta parameters.

### 4.1 A Loss-Based Theoretical Framework

In this section, we view the delta parameter editing operation as a perturbation to the post-trained model. We define the editing perturbation as  $e \triangleq \Delta \widetilde{W}_{edit} - \Delta W$ , which describes the deviation of the edited delta parameters from the original ones. We focus on the loss change caused by editing:

$$\Delta \mathcal{L} \triangleq \mathcal{L}(W_{edit}) - \mathcal{L}(W_{post}) \quad (4)$$

The goal of editing is to control  $|\Delta \mathcal{L}|$  and avoid significant performance degradation.

To analyze how editing affects  $\Delta \mathcal{L}$ , we apply a second-order Taylor expansion:

$$\Delta \mathcal{L} \approx g^\top e + \frac{1}{2} e^\top C e \quad (5)$$

where  $g = \nabla \mathcal{L}(W_{post})$  is the gradient and  $C \succeq 0$  is a positive semi-definite curvature proxy (e.g., Gauss-Newton or Fisher information; we use this instead of the exact Hessian which may be indefinite in deep networks). This expansion decomposes  $\Delta \mathcal{L}$  into a first-order term  $g^\top e$  and a second-order term  $\frac{1}{2} e^\top C e$ , which we analyze in the following subsections.

### 4.2 Editing Intensity

In this part, we focus on the generalized editing formulation defined in Equation 3. We analyze how the choice of  $(p, k)$  affects the loss change  $\Delta \mathcal{L}$ . The editing perturbation  $e$  can be written as:

$$e_i = \begin{cases} (k-1)\Delta w_i & \text{with probability } p \\ \frac{p(1-k)}{1-p}\Delta w_i & \text{with probability } 1-p \end{cases}$$

For the first-order term  $g^\top e = \sum_i g_i e_i$ , we can compute its expectation and variance (detailed

derivation in Appendix C):

$$\begin{aligned} \mathbb{E}[g^\top e] &= 0 \\ \text{Var}(g^\top e) &= \frac{p}{1-p}(1-k)^2 \cdot \sum_i (g_i \Delta w_i)^2 \end{aligned}$$

The expectation being zero indicates that the rescale operation centers the first-order contribution in expectation. This is the foundation for why this type of editing can largely preserve performance. The variance is non-zero and scales with  $(p, k)$ , which means that as the editing becomes more aggressive, the model after a single editing realization may have larger loss deviation.

For the second-order term  $\frac{1}{2}e^\top C e$ , adopting a diagonal approximation  $e^\top C e \approx \sum_i s_i e_i^2$  where  $s_i \geq 0$ , the expectation is:

$$\mathbb{E}\left[\frac{1}{2}\sum_i s_i e_i^2\right] = \frac{1}{2} \cdot \frac{p}{1-p}(1-k)^2 \cdot \sum_i s_i (\Delta w_i)^2$$

Since  $s_i \geq 0$ , this term is always non-negative, representing a curvature cost that accumulates with the perturbation magnitude. The expectation scales with  $(p, k)$  through the factor  $\frac{p}{1-p}(1-k)^2$ , and with the model/task through  $\sum_i s_i (\Delta w_i)^2$ .

Both the variance of the first-order term and the expectation of the second-order term share a common factor that depends on  $(p, k)$ . This factor directly controls the magnitude of loss change: larger values lead to larger variance in the first-order term and larger expected cost in the second-order term. We define the **editing intensity**:

$$\mathcal{I}(p, k) \triangleq \frac{p}{1-p}(1-k)^2 \quad (6)$$

For a fixed model and task,  $\text{Var}(g^\top e) \propto \mathcal{I}$  and  $\mathbb{E}[e^\top C e] \propto \mathcal{I}$ . Thus, larger  $\mathcal{I}$  leads to larger loss fluctuations and larger expected curvature cost, making  $\Delta\mathcal{L}$  more likely to increase, and consequently causing the edited model to deviate further from the post-trained model.

The editing intensity explains the boundary phenomena observed in Section 3. First, when  $p \rightarrow 1$ ,  $\mathcal{I} \rightarrow \infty$  due to the  $\frac{p}{1-p}$  factor, which explains why high drop rates lead to instability regardless of the value of  $k$ . Second, when  $k < 0$  (sign-flip),  $(1-k)^2 > 1$ . For example, when  $k = -1$ ,  $(1-k)^2 = 4$ , which is four times larger than when  $k = 0$  (original DARE). This means that at the same drop rate  $p$ , sign-flip has significantly

larger editing intensity than magnitude-only editing, which explains why sign-flip enters the unstable region earlier and can only tolerate smaller values of  $p$ .

To validate the proposed editing intensity, we evaluate it on LLaMA-3-8B-Instruct using the GSM8K benchmark. Specifically, we sweep over a wide range of  $(p, k)$  to generate a large collection of edited models and measure their downstream performance as a function of the corresponding editing intensity. As shown in Figure 7, editing intensity exhibits a clear negative correlation with performance: as  $\mathcal{I}$  increases, performance consistently decreases. Moreover, when  $\mathcal{I}$  remains small, the edited models stay close to the post-trained model in terms of performance. We further show the relationship on a log-scale in Figure 14, where we observe that under this setting, when  $\mathcal{I} \leq 2$ , performance shows nearly no degradation.

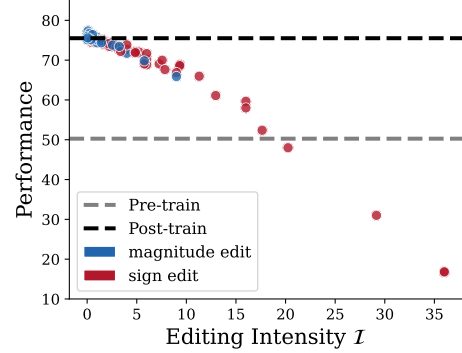


Figure 7: The relationship between editing intensity and performance on LLaMA-3-8B-Instruct using GSM8K benchmark.

### 4.3 Analysis on Magnitude Properties

In this subsection, we analyze how the magnitude properties identified in Section 3 relate to the loss change  $\Delta\mathcal{L}$ . For the first-order term  $g^\top e = \sum_i g_i e_i$ , we decompose it using absolute values and signs:  $g^\top e = \sum_i |g_i| |e_i| a_i$ , where  $a_i = \text{sign}(g_i) \cdot \text{sign}(e_i) \in \{-1, +1\}$ . Defining the normalized weights  $w_i = \frac{|g_i|}{\sum_j |g_j|}$ , we have

$$g^\top e = \left( \sum_i |g_i| \right) \left( \mathbb{E}_w[|e|] \mathbb{E}_w[a] + \text{Cov}_w(|e|, a) \right)$$

Here  $\mathbb{E}_w[|e|]$  measures a gradient-weighted mean magnitude of the perturbation, and  $\mathbb{E}_w[a]$  measures the average sign alignment; the covariance term captures residual heterogeneity. This connects the importance of mean magnitude to the first-order contribution.

For the second-order term, we adopt a diagonal approximation  $e^\top C e \approx \sum_i s_i e_i^2$ , where  $s_i \geq 0$ . This can be decomposed as (detailed derivation in Appendix D):

$$\sum_i s_i e_i^2 = N \bar{s} \bar{e}^2 + N \cdot \text{Cov}(s, e^2),$$

where  $N$  is the number of coordinates,  $\bar{s} = \frac{1}{N} \sum_i s_i$ , and  $\bar{e}^2 = \frac{1}{N} \sum_i e_i^2$ . The first term  $N \bar{s} \bar{e}^2$  captures the global intensity of the perturbation (which increases as the mean magnitude of  $|e|$  increases), and the second term  $N \cdot \text{Cov}(s, e^2)$  reflects the placement of energy across coordinates, which is related to relative relationships.

This analysis explains our experimental findings in Section 3. The mean magnitude affects  $\mathbb{E}_w[|e|]$  in the first-order term and  $\bar{e}^2$  in the second-order term. Relative relationships, when disrupted through operations like shuffling, mainly affect the covariance term  $\text{Cov}(s, e^2)$ , explaining their moderate impact. Specific values and distribution shape, within a reasonable editing range, do not systematically change the mean magnitude or placement structure, thus showing limited impact on  $\Delta \mathcal{L}$ .

## 5 Related Work

**Post-training of Large-Scale Models** Post-training is widely adopted to achieve a pre-trained backbone toward downstream capability and alignment objectives (Dodge et al., 2020; Zhao et al., 2025; Team, 2025). Concretely, the training signal may come from supervised demonstrations (Zhao et al., 2024; Lambert et al., 2025; Moshkov et al., 2025), preference-based optimization, e.g., PPO or DPO, (Ouyang et al., 2022; DeepSeek-AI et al., 2024; Xu et al., 2024; Wang et al., 2024), or verifiable feedback produced by rule-based or model-based verifiers, (Shao et al., 2024; Yu et al., 2025). The effectiveness of post-training can be denoted by the delta parameters, which represent the difference between post-trained and pre-trained parameters (Ilharco et al., 2023; Yu et al., 2024). Given the close correlations between delta parameters and the post-training process, investigating the properties of delta parameters becomes particularly important. In this paper, we discovered the plasticity of delta parameters, suggesting that the effects of post-training can be approximately preserved under diverse parameter configurations.

**Delta Parameter Editing** Delta parameter editing has been explored for various purposes in recent

years. One line of work focuses on model merging, which aims to combine multiple post-trained models into a single model. DARE (Yu et al., 2024) reduces parameter conflicts by randomly dropping delta parameters and rescaling the rest. DELLA-Merging (Deep et al., 2024) extends DARE with magnitude-aware dropping. TIES-Merging (Yadav et al., 2023) resolves sign conflicts and retains only large-magnitude parameters. Twin-Merging (Lu et al., 2024) applies singular value decomposition to extract task-specific knowledge. Another line of work focuses on *model compression*. BitDelta (Liu et al., 2024) quantizes delta parameters to 1-bit by preserving only signs and a shared magnitude scalar. A third line of work focuses on *model enhancement*. EXPO (Zheng et al., 2024) extrapolates delta parameters with a scaling factor to improve alignment performance. While these methods achieve their respective goals through different operations, there remains limited understanding of what properties of delta parameters are essential for maintaining model performance. Our work aims to investigate this question through systematic experiments and provide insights that complement existing methods.

## 6 Conclusion

In this work, we investigated the properties of delta parameters in post-trained models along magnitude and sign. Through experiments across instruct language models, reasoning language models, and vision models, we find that delta parameters exhibit considerable plasticity. In the magnitude dimension, we observe that within a reasonable editing range, the mean is the most sensitive factor, while individual values and distribution shape show less impact. In the sign dimension, we find that a substantial proportion of signs can be flipped while maintaining reasonable performance. To understand these phenomena, we developed a loss-based theoretical framework using second-order Taylor expansion. This framework introduces the concept of editing intensity, which helps explain the stability boundaries of different editing operations, and identifies mean and relative relationships as key factors from a theoretical perspective. Our findings provide insights for understanding and designing delta parameter editing methods.

## 620 Limitations

621 In this work, we systematically investigated the  
622 plasticity of delta parameters across many models.  
623 However, our experiments did not include mixture-  
624 of-experts (MoE) architectures or larger-scale mod-  
625 els (e.g., 70B or above), where the properties of  
626 delta parameters may exhibit different patterns.

## 627 References

628 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,  
629 Henrique Ponde De Oliveira Pinto, Jared Kaplan,  
630 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg  
631 Brockman, and 1 others. 2021. Evaluating large  
632 language models trained on code. *arXiv preprint*  
633 *arXiv:2107.03374*.

634 Shoufa Chen, Chongjian Ge, Zhan Tong, Jiangliu Wang,  
635 Yibing Song, Jue Wang, and Ping Luo. 2022. Adapt-  
636 former: Adapting vision transformers for scalable  
637 visual recognition. In *Advances in Neural Informa-*  
638 *tion Processing Systems 35*.

639 Gong Cheng, Junwei Han, and Xiaoqiang Lu. 2017.  
640 Remote sensing image scene classification: Bench-  
641 mark and state of the art. *Proceedings of the IEEE*,  
642 105(10):1865–1883.

643 Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos,  
644 Sammy Mohamed, and Andrea Vedaldi. 2014. De-  
645 scribing textures in the wild. In *Proceedings of the*  
646 *IEEE conference on computer vision and pattern*  
647 *recognition*, pages 3606–3613.

648 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,  
649 Ashish Sabharwal, Carissa Schoenick, and Oyvind  
650 Tafjord. 2018. Think you have solved question an-  
651 swering? try arc, the ai2 reasoning challenge. *arXiv*  
652 *preprint arXiv:1803.05457*.

653 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
654 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
655 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
656 Nakano, and 1 others. 2021. Training verifiers  
657 to solve math word problems. *arXiv preprint*  
658 *arXiv:2110.14168*.

659 Pala Tej Deep, Rishabh Bhardwaj, and Soujanya Po-  
660 ria. 2024. Della-merging: Reducing interference in  
661 model merging through magnitude-based sampling.  
662 *CoRR*, abs/2406.11617.

663 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,  
664 Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
665 Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,  
666 Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-  
667 hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others.  
668 2025. [Deepseek-r1: Incentivizing reasoning capa-](#)  
669 [bility in llms via reinforcement learning](#). *Preprint*,  
670 *arXiv:2501.12948*.

DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingx-  
uan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng,  
Chong Ruan, Damai Dai, Daya Guo, Dejian Yang,  
Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli  
Luo, Guangbo Hao, Guanting Chen, and 138 others.  
2024. [Deepseek-v2: A strong, economical, and effi-](#)  
[cient mixture-of-experts language model](#). *Preprint*,  
*arXiv:2405.04434*.

Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali  
Farhadi, Hannaneh Hajishirzi, and Noah A. Smith.  
2020. Fine-tuning pretrained language models:  
Weight initializations, data orders, and early stop-  
ping. *CoRR*, abs/2002.06305.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,  
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,  
Akhil Mathur, Alan Schelten, Amy Yang, Angela  
Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,  
Archi Mitra, Archie Sravankumar, Artem Korenev,  
Arthur Hinsvark, Arun Rao, Aston Zhang, and 82  
others. 2024. The llama 3 herd of models. *CoRR*,  
abs/2407.21783.

Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff,  
Dan Jurafsky, and Douwe Kiela. 2024. KTO:  
model alignment as prospect theoretic optimization.  
In *International Conference on Machine Learning*.  
PMLR.

Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and  
Sai Qian Zhang. 2024. Parameter-efficient fine-  
tuning for large models: A comprehensive survey.  
*CoRR*, abs/2403.14608.

Patrick Helber, Benjamin Bischke, Andreas Dengel,  
and Damian Borth. 2019. Eurosat: A novel dataset  
and deep learning benchmark for land use and land  
cover classification. *IEEE Journal of Selected Topics*  
*in Applied Earth Observations and Remote Sensing*,  
12(7):2217–2226.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,  
Mantas Mazeika, Dawn Song, and Jacob Steinhardt.  
2020. Measuring massive multitask language under-  
standing. *arXiv preprint arXiv:2009.03300*.

Gabriel Ilharco, Marco Túlio Ribeiro, Mitchell Worts-  
man, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali  
Farhadi. 2023. Editing models with task arithmetic.  
In *The Eleventh International Conference on Learn-*  
*ing Representations*. OpenReview.net.

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia  
Yan, Tianjun Zhang, Sida Wang, Armando Solar-  
Lezama, Koushik Sen, and Ion Stoica. 2024. Live-  
codebench: Holistic and contamination free eval-  
uation of large language models for code. *arXiv*  
*preprint arXiv:2403.07974*.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-  
sch, Chris Bamford, Devendra Singh Chaplot, Diego  
de Las Casas, Florian Bressand, Gianna Lengyel,  
Guillaume Lample, Lucile Saulnier, Léo Re-  
nard Lavaud, Marie-Anne Lachaux, Pierre Stock,

727	Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. <i>CoRR</i> , abs/2310.06825.	781
728		782
729		783
730	Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 2013. 3d object representations for fine-grained categorization. In <i>Proceedings of the IEEE international conference on computer vision workshops</i> , pages 554–561.	784
731		785
732		786
733		787
734		788
735	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, and 4 others. 2025. <a href="#">Tulu 3: Pushing frontiers in open language model post-training</a> . <i>Preprint</i> , arXiv:2411.15124.	789
736		790
737		791
738		792
739		793
740		794
741		795
742		796
743		797
744	Yann LeCun, Corinna Cortes, and CJ Burges. 2010. Mnist handwritten digit database. <i>ATT Labs [Online]</i> . Available: <a href="http://yann.lecun.com/exdb/mnist">http://yann.lecun.com/exdb/mnist</a> , 2.	798
745		799
746		800
747	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. <a href="#">Let’s verify step by step</a> . <i>Preprint</i> , arXiv:2305.20050.	801
748		802
749		803
750		804
751		805
752	Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. <i>arXiv preprint arXiv:2109.07958</i> .	806
753		807
754		808
755	James Liu, Guangxuan Xiao, Kai Li, Jason D. Lee, Song Han, Tri Dao, and Tianle Cai. 2024. Bitdelta: Your fine-tune may only be worth one bit. <i>CoRR</i> , abs/2402.10193.	809
756		810
757		811
758		812
759	Zhenyi Lu, Chenghao Fan, Wei Wei, Xiaoye Qu, Danyang Chen, and Yu Cheng. 2024. Twin-merging: Dynamic integration of modular expertise in model merging. <i>CoRR</i> , abs/2406.15479.	813
760		814
761		815
762		816
763	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. <i>CoRR</i> , abs/2308.09583.	817
764		818
765		819
766		820
767		821
768		822
769	Ivan Moshkov, Darragh Hanley, Ivan Sorokin, Shubham Toshniwal, Christof Henkel, Benedikt Schifferer, Wei Du, and Igor Gitman. 2025. <a href="#">Aimo-2 winning solution: Building state-of-the-art mathematical reasoning models with openmathreasoning dataset</a> . <i>Preprint</i> , arXiv:2504.16891.	823
770		824
771		825
772		826
773		827
774		828
775	Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Baolin Wu, Andrew Y Ng, and 1 others. 2011. Reading digits in natural images with unsupervised feature learning. In <i>NIPS workshop on deep learning and unsupervised feature learning</i> , volume 2011, page 4. Granada.	829
776		830
777		831
778		832
779		833
780		834
		835
	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> . <i>Preprint</i> , arXiv:2203.02155.	
	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, and 25 others. 2025. <a href="#">Qwen2.5 technical report</a> . <i>Preprint</i> , arXiv:2412.15115.	
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In <i>Advances in Neural Information Processing Systems 36</i> .	
	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2024. Gpqa: A graduate-level google-proof q&a benchmark. In <i>First Conference on Language Modeling</i> .	
	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106.	
	Mark Sandler, Andrey Zhmoginov, Max Vladymyrov, and Andrew Jackson. 2022. Fine-tuning image transformers using learnable memory. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 12145–12154. IEEE.	
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. <a href="#">Deepseekmath: Pushing the limits of mathematical reasoning in open language models</a> . <i>Preprint</i> , arXiv:2402.03300.	
	Johannes Stallkamp, Marc Schlipf, Jan Salmen, and Christian Igel. 2011. The german traffic sign recognition benchmark: a multi-class classification competition. In <i>The 2011 international joint conference on neural networks</i> , pages 1453–1460. IEEE.	
	Qwen Team. 2025. <a href="#">Qwen3 technical report</a> . <i>Preprint</i> , arXiv:2505.09388.	

836	Yuxuan Tong, Xiwen Zhang, Rui Wang, Ruidong Wu, and Junxian He. 2024. Dart-math: Difficulty-aware rejection tuning for mathematical problem-solving. <i>CoRR</i> , abs/2407.13690.	892
837		893
838		894
839		895
840	Junqiao Wang, Zeng Zhang, Yangfan He, Zihao Zhang, Xinyuan Song, Yuyang Song, Tianyu Shi, Yuchen Li, Hengyuan Xu, Kunyu Wu, Xin Yi, Zhongwei Wan, Xinhang Yuan, Zijun Wang, Kuan Lu, Menghao Huo, Tang Jingqun, Guangwu Qian, Keqin Li, and 2 others. 2025. <a href="#">Enhancing code llms with reinforcement learning in code generation: A survey</a> . <i>Preprint</i> , arXiv:2412.20367.	896
841		897
842		898
843		899
844		900
845		901
846		902
847		
848	Zhichao Wang, Bin Bi, Shiva Kumar Pentyala, Kiran Ramnath, Sougata Chaudhuri, Shubham Mehrotra, Zixu, Zhu, Xiang-Bo Mao, Sitaram Asur, Na, and Cheng. 2024. <a href="#">A comprehensive survey of llm alignment techniques: Rlhf, rlaif, ppo, dpo and more</a> . <i>Preprint</i> , arXiv:2407.16216.	903
849		904
850		905
851		906
852		907
853		908
854	Jianxiong Xiao, Krista A Ehinger, James Hays, Antonio Torralba, and Aude Oliva. 2016. Sun database: Exploring a large collection of scene categories. <i>International Journal of Computer Vision</i> , 119:3–22.	910
855		911
856		912
857		913
858	Yi Xin, Siqi Luo, Haodi Zhou, Junlong Du, Xiaohong Liu, Yue Fan, Qing Li, and Yuntao Du. 2024. Parameter-efficient fine-tuning for pre-trained vision models: A survey. <i>CoRR</i> , abs/2402.02242.	914
859		
860		
861		
862	Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, and Yi Wu. 2024. <a href="#">Is DPO superior to PPO for LLM alignment? a comprehensive study</a> . In <i>Forty-first International Conference on Machine Learning</i> .	915
863		916
864		917
865		
866		
867	Prateek Yadav, Derek Tam, Leshem Choshen, Colin A. Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models. In <i>Advances in Neural Information Processing Systems 36</i> .	918
868		919
869		920
870		921
871		
872	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024. <a href="#">Qwen2 technical report</a> . <i>Preprint</i> , arXiv:2407.10671.	
873		
874		
875		
876		
877		
878		
879	Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024. Language models are super mario: Absorbing abilities from homologous models as a free lunch. In <i>International Conference on Machine Learning</i> . PMLR.	
880		
881		
882		
883		
884	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, and 16 others. 2025. <a href="#">Dapo: An open-source llm reinforcement learning system at scale</a> . <i>Preprint</i> , arXiv:2503.14476.	
885		
886		
887		
888		
889		
890		
891		
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i> .	
	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, and 3 others. 2023. A survey of large language models. <i>CoRR</i> , abs/2303.18223.	
	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, and 3 others. 2025. <a href="#">A survey of large language models</a> . <i>Preprint</i> , arXiv:2303.18223.	
	Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. <a href="#">Wildchat: 1m chatGPT interaction logs in the wild</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	
	Chujie Zheng, Ziqi Wang, Heng Ji, Minlie Huang, and Nanyun Peng. 2024. Weak-to-strong extrapolation expedites alignment. <i>CoRR</i> , abs/2404.16792.	
	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. <a href="#">Instruction-following evaluation for large language models</a> . <i>Preprint</i> , arXiv:2311.07911.	

## A Deriving the Rescale Factor in Equation 3

We derive the rescale factor used in Equation 3 from the same expected-output preservation intuition as in DARE. Consider a linear transformation  $h = Wx + b$  with delta parameters  $(\Delta W, \Delta b)$ . We apply a coordinate-wise random scaling to the delta parameters: with probability  $p$  we multiply by  $k$ , and with probability  $1 - p$  we multiply by an unknown factor  $\gamma$ . Let  $M \sim \text{Bernoulli}(p)$  be the binary mask (element-wise) indicating which coordinates take the factor  $k$ .

After editing, the output becomes

$$\hat{h} = (W_{\text{pre}} + \Delta\tilde{W})x + (b_{\text{pre}} + \Delta\tilde{b}),$$

where,

$$\Delta\tilde{W} = k \cdot M \odot \Delta W + \gamma \cdot (1 - M) \odot \Delta W,$$

$$\Delta\tilde{b} = k \cdot M \odot \Delta b + \gamma \cdot (1 - M) \odot \Delta b.$$

Taking expectation over the editing randomness (i.e., over  $M$ ), we have  $\mathbb{E}[M] = p$  and  $\mathbb{E}[1 - M] = 1 - p$ . Thus

$$\mathbb{E}[\Delta\tilde{W}] = (pk + (1 - p)\gamma)\Delta W,$$

$$\mathbb{E}[\Delta\tilde{b}] = (pk + (1 - p)\gamma)\Delta b.$$

DARE-style expected-output preservation requires the expected delta contribution to match the original delta contribution, i.e.,  $\mathbb{E}[\Delta\tilde{W}] = \Delta W$  and  $\mathbb{E}[\Delta\tilde{b}] = \Delta b$ . This yields a single scalar condition:

$$pk + (1 - p)\gamma = 1.$$

Solving for  $\gamma$  gives the rescale factor used in Equation 3:

$$\gamma = \frac{1 - kp}{1 - p}.$$

## B Full Experimental Results

We conduct a thorough experimental validation on the plasticity of delta parameters. The results of LLaMA3-8B-Instruct, Mistral-7B-Instruct-v0.3, ViT-B-32 and Qwen3-1.7B across eight benchmarks are presented in Figure 8, Figure 9, Figure 10, and Figure 11, Figure 12, respectively.

The full results of LLaMA-3-8B-Instruct on all datasets for experiments with varying magnitude mean and distribution shape are shown in Figure 13.

## C Derivations for Editing Intensity

This appendix provides derivations for the results in Section 4. We start from the two-point distribution of the editing perturbation already given in the main text:

$$e_i = \begin{cases} (k - 1)\Delta w_i & \text{with probability } p, \\ \frac{p(1-k)}{1-p}\Delta w_i & \text{with probability } 1 - p. \end{cases}$$

We take expectation/variance over the editing randomness, treating  $(g, \Delta W)$  as fixed for the local surrogate.

### C.1 First-order term: $\mathbb{E}[g^\top e] = 0$

We compute  $\mathbb{E}[e_i]$  directly:

$$\mathbb{E}[e_i] = p(k - 1)\Delta w_i + (1 - p)\frac{p(1 - k)}{1 - p}\Delta w_i$$

$$= p(k - 1)\Delta w_i + p(1 - k)\Delta w_i$$

$$= 0.$$

Therefore,

$$\mathbb{E}[g^\top e] = \mathbb{E}\left[\sum_i g_i e_i\right] = \sum_i g_i \mathbb{E}[e_i] = 0.$$

### C.2 Variance of the first-order term

We assume the coordinate-wise editing randomness is independent across  $i$  (e.g., induced by an i.i.d. Bernoulli mask). Since  $\mathbb{E}[e_i] = 0$ , we have  $\text{Var}(e_i) = \mathbb{E}[e_i^2]$ . First compute  $\mathbb{E}[e_i^2]$ :

$$\mathbb{E}[e_i^2] = p(k - 1)^2(\Delta w_i)^2 + (1 - p)\left(\frac{p(1 - k)}{1 - p}\right)^2(\Delta w_i)^2$$

$$= (1 - k)^2(\Delta w_i)^2\left(p + \frac{p^2}{1 - p}\right)$$

$$= (1 - k)^2(\Delta w_i)^2 \cdot \frac{p}{1 - p}.$$

Now let  $X_i \triangleq g_i e_i$ . Under independence across  $i$ ,  $\text{Var}(\sum_i X_i) = \sum_i \text{Var}(X_i)$ . Thus:

$$\text{Var}(g^\top e) = \text{Var}\left(\sum_i g_i e_i\right) = \sum_i \text{Var}(g_i e_i)$$

$$= \sum_i g_i^2 \text{Var}(e_i) = \sum_i g_i^2 \mathbb{E}[e_i^2]$$

$$= \frac{p}{1 - p}(1 - k)^2 \sum_i (g_i \Delta w_i)^2.$$

### C.3 Expected diagonal second-order term and the common factor

Under the diagonal curvature surrogate used in the main text,  $e^\top C e \approx \sum_i s_i e_i^2$  with  $s_i \geq 0$ . Taking expectation:

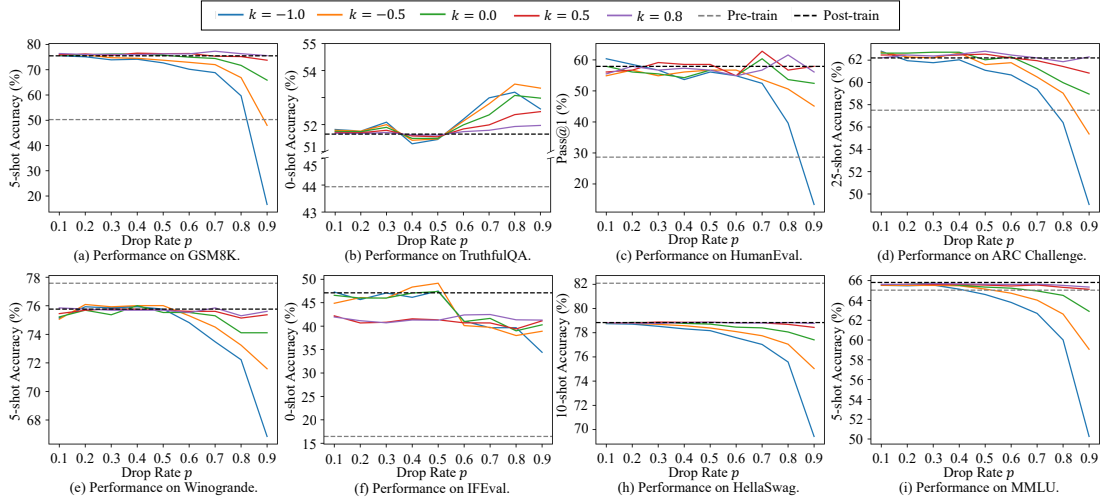


Figure 8: The performance of LLaMA3-8B-Instruct on the all benchmarks under varying  $p$  and  $k$ .

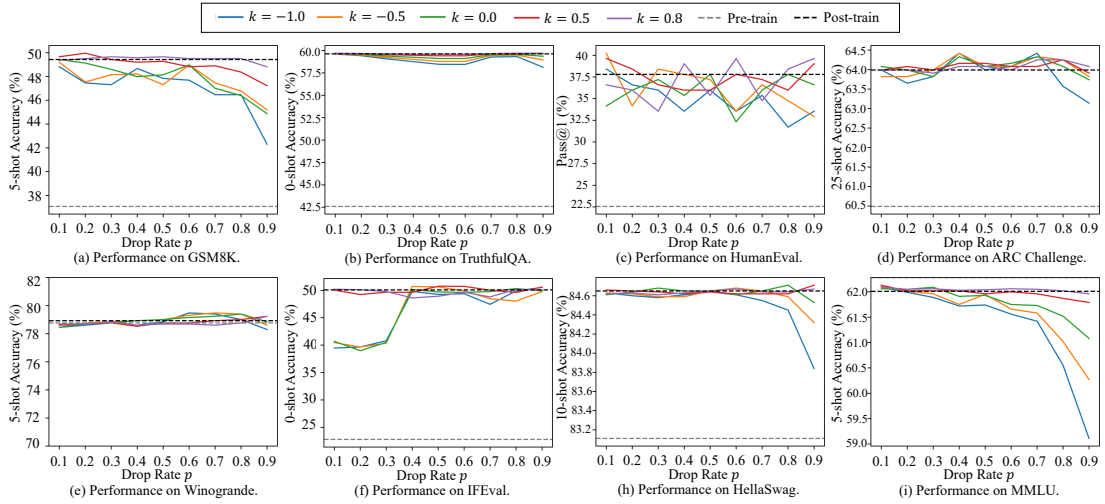


Figure 9: The performance of Mistral-7B-Instruct-v0.3 on the all benchmarks under varying  $p$  and  $k$ .

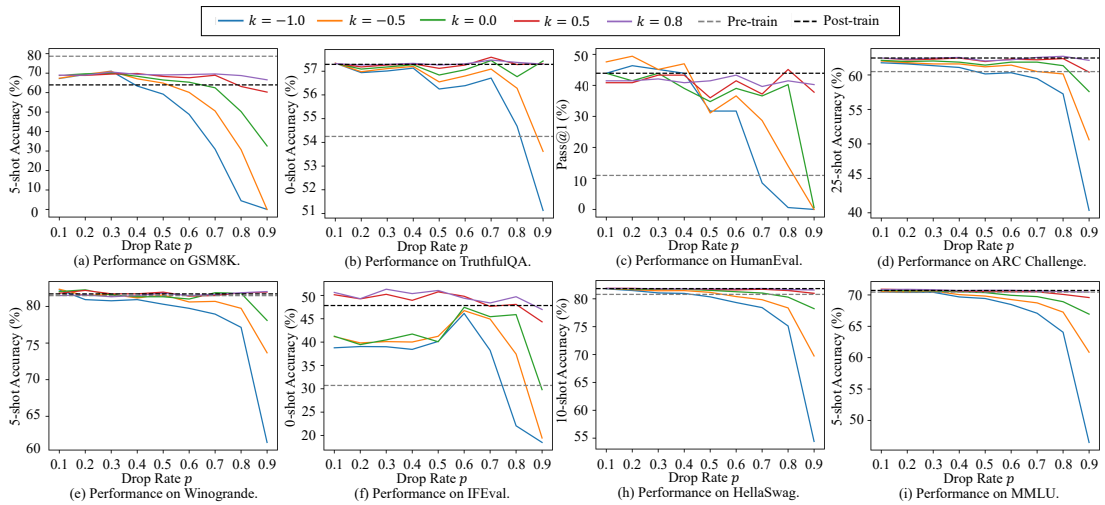


Figure 10: The performance of Qwen2-7B-Instruct on the all benchmarks under varying  $p$  and  $k$ .

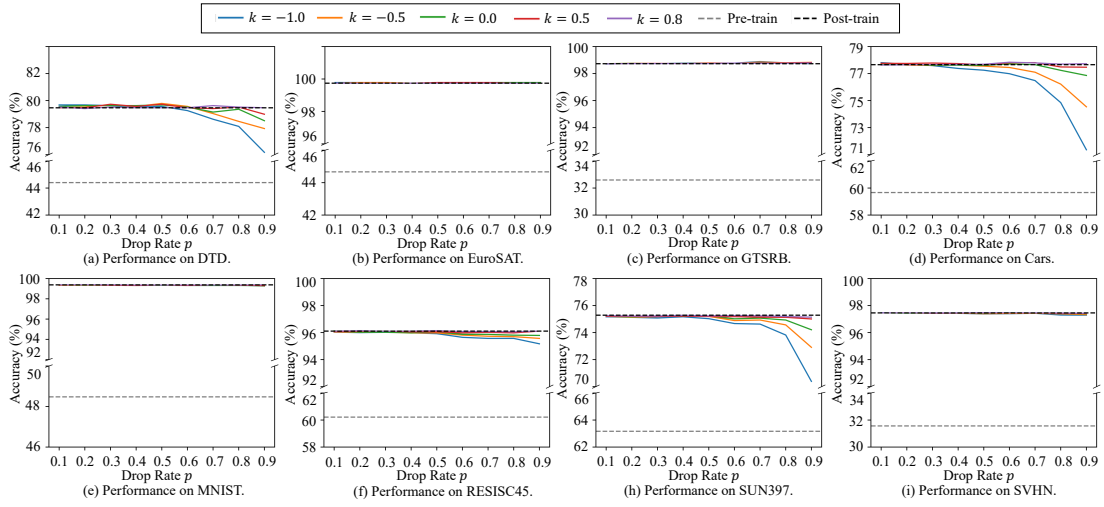


Figure 11: The performance of ViT-B-32 on the all benchmarks under varying  $p$  and  $k$ .

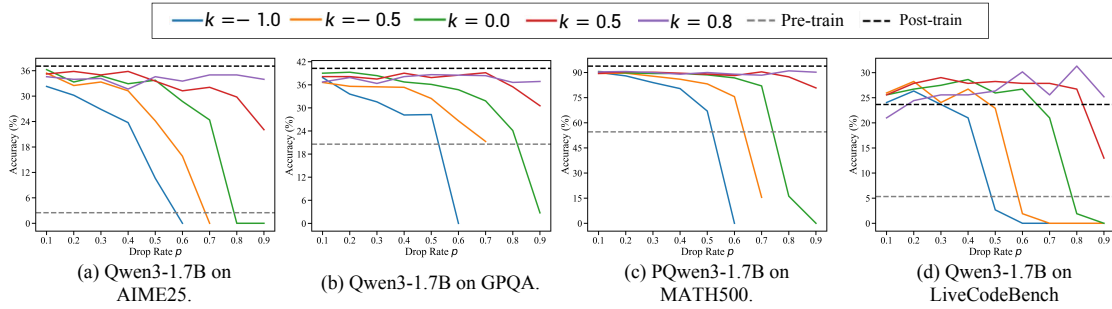


Figure 12: The performance of Qwen3-1.7B on the all benchmarks under varying  $p$  and  $k$ .

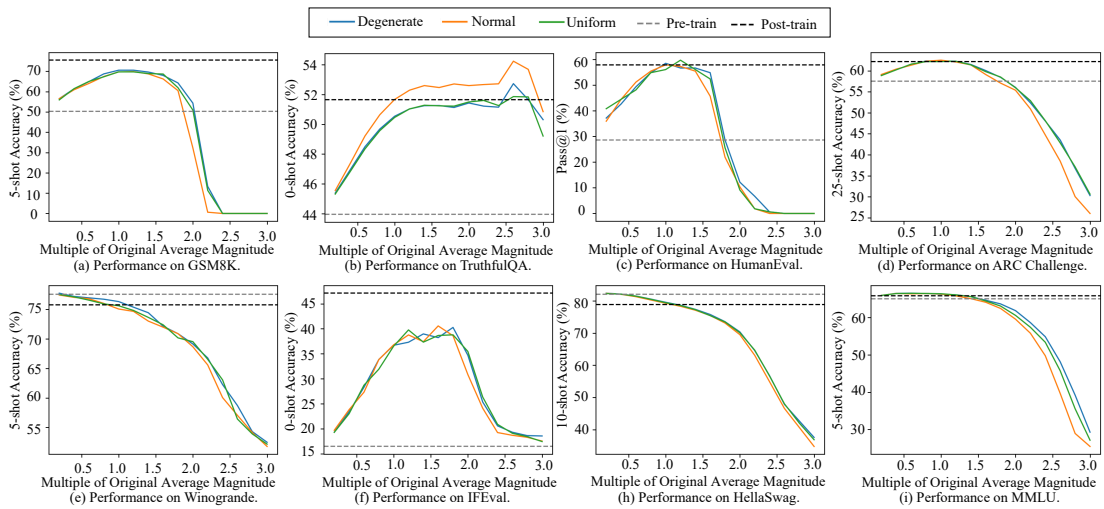


Figure 13: The full results on all datasets for experiments with varying magnitude mean and distribution shape.

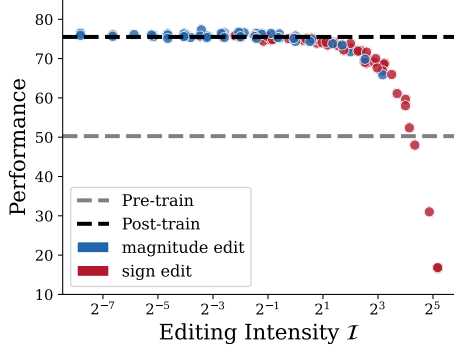


Figure 14: The relationship between editing intensity and performance in log-scale on LLaMA-3-8B-Instruct using GSM8K benchmark.

$$\begin{aligned} \mathbb{E}\left[\frac{1}{2}\sum_i s_i e_i^2\right] &= \frac{1}{2}\sum_i s_i \mathbb{E}[e_i^2] \\ &= \frac{1}{2} \cdot \frac{p}{1-p}(1-k)^2 \sum_i s_i (\Delta w_i)^2. \end{aligned}$$

This shows that both  $\text{Var}(g^\top e)$  and  $\mathbb{E}[\sum_i s_i e_i^2]$  share the same  $(p, k)$ -dependent factor. We therefore define the editing intensity:

$$\mathcal{I}(p, k) \triangleq \frac{p}{1-p}(1-k)^2.$$

## D Decomposition Identities

### D.1 Weighted covariance form for the first-order term

Define  $a_i = \text{sign}(g_i)\text{sign}(e_i) \in \{-1, +1\}$  so that  $g_i e_i = |g_i||e_i|a_i$ . Let  $w_i = \frac{|g_i|}{\sum_j |g_j|}$  with  $\sum_i w_i = 1$ . Then:

$$\begin{aligned} g^\top e &= \sum_i |g_i||e_i|a_i \\ &= \left(\sum_i |g_i|\right) \sum_i w_i |e_i|a_i \\ &= \left(\sum_i |g_i|\right) \mathbb{E}_w[|e|a]. \end{aligned}$$

Using the weighted covariance identity  $\mathbb{E}_w[XY] = \mathbb{E}_w[X]\mathbb{E}_w[Y] + \text{Cov}_w(X, Y)$  with  $X = |e|$  and  $Y = a$  gives the expression in the main text.

### D.2 Unweighted covariance form for the diagonal second-order term

We use the unweighted (coordinate-wise) averages  $\bar{s} = \frac{1}{N}\sum_i s_i$  and  $\bar{e^2} = \frac{1}{N}\sum_i e_i^2$ . Define the unweighted covariance as

$$\text{Cov}(s, e^2) \triangleq \frac{1}{N}\sum_i (s_i - \bar{s})(e_i^2 - \bar{e^2}). \quad 1012$$

Then: 1013

$$\frac{1}{N}\sum_i s_i e_i^2 = \bar{s}\bar{e^2} + \text{Cov}(s, e^2), \quad 1014$$

$$\sum_i s_i e_i^2 = N\bar{s}\bar{e^2} + N\text{Cov}(s, e^2),$$

which is the decomposition used in Section 4. 1015

## E The Use of Large Language Models 1016

We utilized LLMs to aid and polish writing. 1017