
Interaction-Aware Influence Functions for Group Attribution

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Influence functions approximate how removing a training example changes a
2 quantity of interest, called the target function, such as a held-out loss. To estimate
3 the influence of a group of examples, the standard practice is to sum the individual
4 influences of its members. However, this sum does not capture how examples
5 jointly affect the target: a pair of examples may be redundant or complementary,
6 but the sum cannot distinguish these cases. We propose an interaction-aware
7 influence function that characterizes how interactions between examples influence
8 the target. By expanding the target to second order around the trained parameters,
9 we obtain an estimator that augments the standard sum with a pairwise interaction
10 term that captures the alignment between two examples' effects on the target. We
11 empirically evaluate our estimator in two settings. First, on six dataset-model
12 pairs spanning logistic regression, MLPs, and ResNet-9, our estimator tracks
13 leave-group-out retraining substantially better than first-order influence across
14 all settings. Second, when used as a greedy selection rule for instruction-tuning
15 data on Llama-3.1-8B, it beats prior influence-based and representation-similarity
16 baselines on five of seven downstream tasks, in a regime where standard influence-
17 based selection underperforms random selection. Code is available at https://anonymous.4open.science/r/Interaction_IF-45D6.
18

19 1 Introduction

20 Quantifying how training examples shape a trained model is a foundational problem in machine
21 learning, with applications including selecting valuable training data [18, 59, 64], understanding
22 what kinds of data benefit a model [22, 24], and fairly compensating data providers [10, 20]. The
23 conceptual gold standard is leave-one-out retraining: remove a training example, retrain the model,
24 and measure how some quantity of interest has changed. Since retraining for every example is
25 prohibitive, influence functions [29] estimate this counterfactual in closed form by modeling two
26 things in turn: how reweighting a training example would shift the trained parameters, and how those
27 shifted parameters would change the quantity of interest. The quantity of interest, called the *target*
28 *function*, is left to the user. It can be the loss on a held-out example [29], an average over a validation
29 set, the likelihood of a downstream task answer [10, 59], or any other smooth function of the model.
30 This flexibility is what lets a single approximation serve so many applications.

31 Most of these applications, including subset selection, group-robustness analysis, and data valuation,
32 ultimately ask about *groups* of examples rather than single ones. The textbook extension scores a
33 group as the sum of its members' individual influences [30], implicitly assuming examples contribute
34 independently. This assumption is rarely true. Two near-duplicate examples each look highly
35 influential on their own, yet adding both has roughly the same effect as adding one, so additive
36 scoring double-counts the overlap and overestimates the group's true effect [25, 26]. The same blind
37 spot makes top-*k* influence-based data selection produce subsets that are highly redundant rather than

38 collectively useful [67], since the most individually influential examples tend to be similar. Basu et al.
 39 [5] previously addressed this with a second-order correction that makes the estimator non-additive,
 40 but their construction refines only how the trained parameters shift, missing how the target function
 41 responds to those shifts.

42 We argue that group attribution is fundamentally about how training examples affect the target.
 43 Therefore, interactions between examples should be characterized in terms of their joint effect on
 44 the target, not just on the parameter shift. To this end, we propose an *interaction-aware influence*
 45 *function* that augments the standard sum of individual influences with a pairwise interaction term over
 46 examples in the group. The pairwise term captures the joint effect on the target: it is positive when the
 47 two examples push the model in similar directions and negative when they push in opposite directions.
 48 We further derive a closed-form factorization of the interaction term across classification settings,
 49 proving that two similar examples are jointly complementary when their labels differ and redundant
 50 when their labels agree, formally recovering classical principles in supervised learning [7, 15, 39, 46].
 51 As an application, we use the interaction term to discount candidates whose effects overlap with
 52 already selected examples, addressing the redundancy in top- k influence-based selection.

53 We evaluate the method in two settings. First, we test whether our estimator faithfully tracks
 54 leave-group-out retraining across six dataset–model pairs, by assessing Spearman rank correlation
 55 with ground-truth retraining effects. It substantially outperforms prior influence-based training data
 56 attribution methods, demonstrating the importance of accounting for interactions in group attribution.
 57 Second, we evaluate the data selection method at scales ranging from small MLPs to instruction-
 58 tuning of Llama-3.1-8B [21]. On small-scale models, our method selects subsets with class diversity
 59 comparable to random selection, whereas baselines collapse onto a few classes and fall below random.
 60 In instruction-tuning data selection, our method consistently improves over random selection and
 61 achieves the strongest overall performance on five of seven downstream tasks.

62 2 Preliminary

63 We first review the classical influence function and its standard extension to group attribution. We
 64 summarize the notation used throughout the paper in Table 2 of Appendix A. Let $\mathcal{D} = \{z_i\}_{i=1}^N$ be
 65 the training set of size N , where each example $z_i = (x_i, y_i)$ consists of an input x_i and a label y_i ,
 66 and let $\ell(z_i, \theta)$ denote the loss of example z_i at parameter $\theta \in \Theta$. For any subset $A \subseteq \mathcal{D}$, define the
 67 empirical risk and its minimizer as

$$\mathcal{R}_A(\theta) := \frac{1}{|A|} \sum_{z_i \in A} \ell(z_i, \theta), \quad \hat{\theta}_A := \arg \min_{\theta} \mathcal{R}_A(\theta). \quad (1)$$

68 We write $\hat{\theta} := \hat{\theta}_{\mathcal{D}}$ for the parameter trained on the full dataset. For a subset $S \subseteq \mathcal{D}$, the parameter
 69 obtained by retraining after removing S is $\hat{\theta}_{\mathcal{D} \setminus S}$.

70 Let $f : \Theta \rightarrow \mathbb{R}$ be a differentiable scalar target function of the trained parameters that we treat as a
 71 quantity to be minimized, such as the loss on a test example or the average loss on a validation set.¹
 72 The exact leave-group-out effect of removing S is

$$\mathcal{I}^-(S) := f(\hat{\theta}_{\mathcal{D} \setminus S}) - f(\hat{\theta}). \quad (2)$$

73 For a singleton $S = \{z_i\}$, this reduces to the usual leave-one-out effect. Computing $\mathcal{I}^-(S)$ exactly
 74 requires retraining the model without S , which is expensive even for a single example and infeasible
 75 for many candidate groups.

76 Influence functions approximate this retraining effect through infinitesimal reweighting [29]. Instead
 77 of removing S directly, consider the parameter obtained after adding a small weight ϵ to the losses of
 78 examples in S :

$$\hat{\theta}(\epsilon; S) := \arg \min_{\theta} \left\{ \frac{1}{N} \sum_{z_i \in \mathcal{D}} \ell(z_i, \theta) + \epsilon \sum_{z_i \in S} \ell(z_i, \theta) \right\}. \quad (3)$$

79 At $\epsilon = 0$, this recovers the original parameter $\hat{\theta}$. Since each example has weight $1/N$ in the empirical
 80 risk, setting $\epsilon = -1/N$ removes the contribution of every example in S . Thus, leave-group-out
 81 retraining can be viewed as a finite step along the reweighting path $\hat{\theta}(\epsilon; S)$.

¹A target to be maximized can be handled by negating f .

82 The influence function replaces this finite step with two first-order approximations. Let $H :=$
83 $\frac{1}{N} \sum_{z_i \in \mathcal{D}} \nabla_{\hat{\theta}}^2 \ell(z_i, \hat{\theta})$ be the Hessian of the training objective. Assuming that $\ell(z_i, \cdot)$ is twice continu-
84 ously differentiable for every $z_i \in \mathcal{D}$ and that H is nonsingular at $\hat{\theta}$, the first approximation linearizes
85 the reweighting path $\hat{\theta}(\epsilon; S)$ in ϵ at $\epsilon = 0$, yielding the parameter shift induced by removing S as

$$\hat{\theta}_{\mathcal{D} \setminus S} - \hat{\theta} \approx \frac{1}{N} H^{-1} \sum_{z_i \in S} \nabla_{\theta} \ell(z_i, \hat{\theta}). \quad (4)$$

86 To simplify the notation, let $u_i := H^{-1} \nabla_{\theta} \ell(z_i, \hat{\theta})$ denote the per-example parameter shift induced
87 by z_i , and let $u_S := \sum_{z_i \in S} u_i$ denote its aggregate over S . The second approximation linearizes the
88 target function f around $\hat{\theta}$, giving the standard group influence-function estimate of $\mathcal{I}^-(S)$:

$$\hat{\mathcal{I}}_{\text{lin}}^-(S) := \frac{1}{N} \nabla_{\theta} f(\hat{\theta})^{\top} u_S. \quad (5)$$

89 The derivations of Equation (4) and Equation (5) are provided in Appendix B. Letting $\hat{\mathcal{I}}_{\text{lin}}^-(z_i)$ denote
90 the singleton case, the standard group estimate is additive: $\hat{\mathcal{I}}_{\text{lin}}^-(S) = \sum_{z_i \in S} \hat{\mathcal{I}}_{\text{lin}}^-(z_i)$. This additivity
91 is a consequence of the first-order approximation, not a property of the exact retraining effect $\mathcal{I}^-(S)$.
92 As a result, standard influence functions cannot capture interactions among examples in a group,
93 which motivates the interaction-aware approximation developed in the next section.

94 3 Method

95 We now develop our interaction-aware influence function, which captures how examples in a group
96 jointly affect a target function. Section 3.1 derives a group-attribution estimator by expanding the
97 target function to second order, yielding the standard additive influence plus a pairwise interaction
98 term. Section 3.2 then extends the estimator to a greedy data-selection procedure that discounts
99 candidates whose effect overlaps with already-selected examples. Full derivations and proofs for this
100 section are provided in Appendix C.

101 3.1 Interaction-aware influence functions

102 The standard influence function loses interaction information because the first-order approximation
103 treats each example’s effect as independent of the others. We extend the first-order approximation to
104 a second-order one to recover the interaction information.

105 **Second-order expansion of the target function.** We expand the target function f to second order
106 around $\hat{\theta}$:

$$f(\hat{\theta}_{\mathcal{D} \setminus S}) - f(\hat{\theta}) = \nabla_{\theta} f(\hat{\theta})^{\top} \delta_S + \frac{1}{2} \delta_S^{\top} H_f \delta_S + O(\|\delta_S\|^3), \quad (6)$$

107 where $\delta_S := \hat{\theta}_{\mathcal{D} \setminus S} - \hat{\theta}$ is the parameter shift induced by removing S and $H_f := \nabla_{\theta}^2 f(\hat{\theta})$ is the
108 Hessian of the target at $\hat{\theta}$. Substituting Equation (4) into Equation (6) yields our interaction-aware
109 influence function, an estimator of $\mathcal{I}^-(S)$:

$$\hat{\mathcal{I}}^-(S) := \underbrace{\frac{1}{N} \nabla_{\theta} f(\hat{\theta})^{\top} u_S}_{\text{first-order influence}} + \underbrace{\frac{1}{2N^2} u_S^{\top} H_f u_S}_{\text{interaction term (ours)}}. \quad (7)$$

110 The first term recovers the standard first-order estimate $\hat{\mathcal{I}}_{\text{lin}}^-(S)$ from Equation (5). The second term,
111 which we call the interaction term, is a curvature correction induced by the target function and is the
112 source of non-additivity across examples in S .

113 The same expansion applies to the data-addition setting, which arises in problems such as data
114 selection and active learning. Adding S to \mathcal{D} flips the sign of the linear term while leaving the
115 quadratic term unchanged, giving

$$\hat{\mathcal{I}}^+(S) := -\frac{1}{N} \nabla_{\theta} f(\hat{\theta})^{\top} u_S + \frac{1}{2N^2} u_S^{\top} H_f u_S. \quad (8)$$

116 The full derivation for the data-addition setting is provided in Appendix C.2. In practice, we use a
117 damped Gauss–Newton Hessian approximation to ensure invertibility, and apply EK-FAC [19, 22]
118 and low-rank gradient projection [10] for scalable Hessian computation. Details of these practical
119 approximations are provided in Appendix D.

120 **Interpreting the pairwise structure of the interaction term.** The interaction term decomposes
 121 into a sum of pairwise contributions between examples in S :

$$u_S^\top H_f u_S = \sum_{z_i, z_j \in S} \kappa(z_i, z_j), \quad \kappa(a, b) := u_a^\top H_f u_b. \quad (9)$$

122 For interpretive clarity, we first consider the case in which H_f is positive definite; $\kappa(a, b)$ is then the
 123 inner product of the parameter shifts u_a and u_b under the metric induced by H_f . A positive value
 124 indicates that the two shifts are aligned under this metric, meaning the examples perturb the parameters
 125 in directions that produce similar changes in the target function, while orthogonal or opposing shifts
 126 produce values near zero or negative. In this sense, $\kappa(a, b)$ quantifies how *similarly* two examples
 127 act on the target. This similarity perspective explains why the additive approximation underestimates
 128 the exact retraining effect on groups of similar examples [25]. Such groups produce large pairwise
 129 interactions that the additive approximation discards, leading to systematic underestimation. In the
 130 more general setting where H_f is indefinite, $\kappa(a, b)$ generalizes from an inner product to a symmetric
 131 bilinear form. We provide a spectral analysis of this case in Appendix C.4.

132 **Relation between the first-order and interaction terms.** A concrete scenario clarifies how the
 133 two terms combine. Suppose we wish to identify a group S whose removal would most reduce a
 134 target loss f . The first-order term, which sums individual influences, picks out examples whose
 135 per-example influence $\frac{1}{N} \nabla_{\theta} f(\hat{\theta})^\top u_i$ is negative, since each such removal individually lowers f . On
 136 a group of near-duplicate examples of this kind, the additive sum predicts a large benefit, scaling
 137 linearly with group size. The interaction term, however, is non-negative whenever H_f is positive
 138 semidefinite, and scales superlinearly with group size when the per-example shifts align. Adding it
 139 back therefore shrinks the predicted benefit, capturing the diminishing returns of removing redundant
 140 examples.

141 In the addition setting, the linear term flips sign while the quadratic term is unchanged. In this
 142 case, the two terms genuinely trade off: the linear term rewards candidates whose addition would
 143 individually lower f , capturing quality, while the interaction term penalizes similarity within the
 144 chosen group, capturing redundancy. This trade-off resonates with a recurring principle in data
 145 selection [67] and active learning [3, 11, 28]: effective subsets pair individually helpful examples
 146 with non-overlapping ones, typically by combining a per-example score with an explicit diversity
 147 term. Here, the same balance emerges directly from a second-order expansion of the target, rather
 148 than from an externally imposed design. The greedy procedure in Section 3.2 exploits this trade-off,
 149 balancing per-example quality against group diversity at every step.

150 **Interpreting the pairwise interaction in binary logistic regression.** To better understand the
 151 pairwise interaction, we analyze a tractable case in which the interaction admits a closed form: binary
 152 logistic regression with ℓ_2 regularization. Let σ denote the sigmoid function, and write $\sigma_i := \sigma(\hat{\theta}^\top x_i)$
 153 for the predicted probability on example i .

154 **Proposition 1** (Closed-form factorization of κ). *For binary logistic regression at $\hat{\theta}$ with $y_i \in \{0, 1\}$,*

$$\kappa(a, b) = (\sigma_a - y_a)(\sigma_b - y_b) \cdot \langle x_a, x_b \rangle_M,$$

155 *where $\langle u, v \rangle_M := u^\top M v$ is the bilinear form induced by $M := H^{-1} H_f H^{-1}$.*

156 The proof relies on the closed-form gradient $\nabla_{\theta} \ell(z_i, \hat{\theta}) = (\sigma_i - y_i) x_i$. Full details are given in
 157 Appendix C.5.

158 Proposition 1 shows that when $\langle x_a, x_b \rangle_M > 0$, so that x_a and x_b are similar under the bilinear
 159 form induced by M , the sign of $\kappa(a, b)$ is determined by the sign of the class-agreement factor
 160 $(\sigma_a - y_a)(\sigma_b - y_b)$. Since $\sigma_i \in [0, 1]$, this factor is positive for same-class pairs and negative for
 161 cross-class pairs. Because the target is a quantity to be minimized, this implies that adding two similar
 162 examples with different labels benefits the classifier beyond their individual effects, while adding
 163 similar examples with the same label yields diminishing returns.

164 This formally recovers two well-known principles of supervised learning: cross-class pairs with
 165 similar features shape the decision boundary and are thus complementary [15, 46, 50], whereas same-
 166 class pairs with similar features provide overlapping evidence and yield diminishing returns [7, 39, 52].
 167 Proposition 1 shows that the interaction term encodes both principles directly, with the sign of $\kappa(a, b)$

Algorithm 1 Greedy selection with interaction-aware influence

Require: Candidate pool $\mathcal{D}_{\text{pool}} = \{z_i\}_{i=1}^P$, budget K , precomputed $\{u_i\}$, $\{w_i := H_f u_i\}$, $\{q_i := u_i^\top w_i\}$, target gradient $\nabla_\theta f(\hat{\theta})$, training set size $N := |\mathcal{D}|$
Ensure: Selected subset S with $|S| \leq K$

- 1: Initialize $S \leftarrow \emptyset$ and accumulator $w \leftarrow \mathbf{0}$, where w tracks $H_f u_S$
- 2: **for** $t = 1, \dots, K$ **do**
- 3: **for** $i \notin S$ **do**
- 4: $m(z_i | S) \leftarrow -\frac{1}{N} \nabla_\theta f(\hat{\theta})^\top u_i + \frac{1}{N^2} w^\top u_i + \frac{1}{2N^2} q_i$
- 5: **end for**
- 6: $i^* \leftarrow \arg \min_{i \notin S} m(z_i | S)$
- 7: $S \leftarrow S \cup \{z_{i^*}\}$, $w \leftarrow w + w_{i^*}$
- 8: **end for**
- 9: **return** S

168 aligning with the role each pair plays in shaping the classifier. For clarity of exposition, we have
169 presented this analysis in the binary logistic regression setting. Both the factorization and these
170 conclusions hold for deep multi-class classifiers, as we develop in Appendix C.6.

171 **Remark** (Comparison with Basu et al. [5]). *Our framework expands the target function to second*
172 *order, while the same second-order treatment can equally be applied to the reweighted minimizer*
173 *$\hat{\theta}(\epsilon; S)$. Basu et al. [5] pursue this alternative for group attribution. The two approaches yield*
174 *structurally distinct estimators: ours decomposes into pairwise contributions between examples in S ,*
175 *which enables the similarity interpretation and the closed-form analysis in Proposition 1, while theirs*
176 *operates at the parameter level without an analogous pairwise structure. The two approaches also*
177 *differ in stability: $\hat{\theta}(\epsilon; S)$ need not be a smooth function of ϵ in non-convex settings, so its second-*
178 *order expansion in ϵ is poorly conditioned, while the target functions we consider, such as the held-out*
179 *loss, are twice continuously differentiable in θ and avoid this source of ill-conditioning. Consistent*
180 *with this analysis, Basu et al. [5] themselves report degraded performance of their estimator on*
181 *non-linear models, a finding that aligns with our empirical comparison in Section 4.*

182 3.2 Interaction-aware data selection

183 We now propose a data selection method that exploits the interaction term to select examples from a
184 candidate pool for addition to the training set. This problem reduces directly to group attribution:
185 finding the group that most reduces the target when added to the training data. Formally, given a
186 model trained on \mathcal{D} and a candidate pool $\mathcal{D}_{\text{pool}}$, we seek a subset $S \subseteq \mathcal{D}_{\text{pool}}$ of size K that minimizes
187 the following objective:

$$S^* := \arg \min_{|S|=K} f(\hat{\theta}_{\mathcal{D} \cup S}), \quad (10)$$

188 where $\hat{\theta}_{\mathcal{D} \cup S}$ denotes the minimizer of the empirical risk on the augmented set $\mathcal{D} \cup S$. Since $f(\hat{\theta})$
189 does not depend on S , subtracting it from the objective preserves the minimizer, so this is equivalent
190 to minimizing the inclusion effect $\mathcal{I}^+(S) := f(\hat{\theta}_{\mathcal{D} \cup S}) - f(\hat{\theta})$. As $\mathcal{I}^+(S)$ is intractable to evaluate
191 exactly, we use its estimator $\hat{\mathcal{I}}^+(S)$ from Equation (8) as our selection criterion.

192 **Greedy procedure.** Exact combinatorial optimization of $\hat{\mathcal{I}}^+(S)$ is intractable, so we instead adopt
193 a greedy procedure that, starting from $S = \emptyset$, repeatedly appends the candidate z_i minimizing
194 $\hat{\mathcal{I}}^+(S \cup \{z_i\})$. We quantify this change through the marginal score of adding $z_i \notin S$ to a current
195 subset S :

$$m(z_i | S) := \hat{\mathcal{I}}^+(S \cup \{z_i\}) - \hat{\mathcal{I}}^+(S). \quad (11)$$

196 Substituting Equation (8) into Equation (11) and using $u_{S \cup \{z_i\}} = u_S + u_i$, which follows from the
197 definition of u_S , to expand the quadratic term gives

$$m(z_i | S) \approx -\frac{1}{N} \nabla_\theta f(\hat{\theta})^\top u_i + \frac{1}{N^2} u_S^\top H_f u_i + \frac{1}{2N^2} u_i^\top H_f u_i. \quad (12)$$

198 The first term is the standard first-order influence of z_i , which captures its individual contribution
199 to reducing the target. The second term is the key novelty of our procedure: it accumulates the

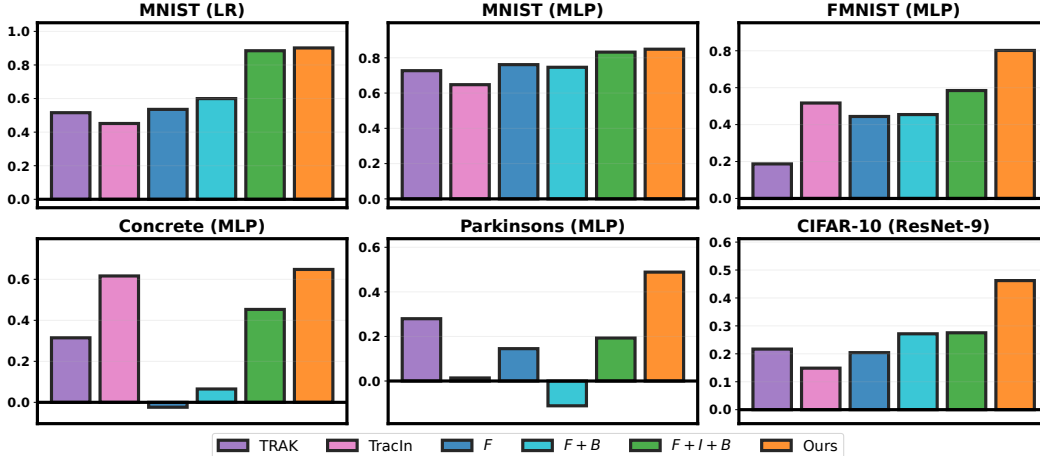


Figure 1: Spearman rank correlation between estimated and ground-truth group influences across six dataset–model pairs, where F , B , and I denote the first-order influence function, the second-order term from Basu et al. [5], and our interaction term, respectively.

200 pairwise interactions analyzed in Section 3.1 between the candidate z_i and each example in the
 201 already-selected subset S . When z_i is similar to examples in S , the second term takes a large positive
 202 value and increases the marginal score. The procedure thereby avoids selecting such candidates. The
 203 third term depends only on the candidate z_i itself, measuring how strongly its own parameter shift
 204 would perturb the target. At each iteration we append the candidate minimizing $m(z_i | S)$ until
 205 $|S| = K$; the full procedure is given in Algorithm 1.

206 The interaction-aware structure of Equation (12) sets our procedure apart from selection based on
 207 the first-order influence [59], which chooses the top- k examples by individual influence score. Since
 208 high-influence examples tend to share characteristics [67], this top- k rule yields redundant subsets and
 209 consequently diminishing returns. Our second term addresses this limitation by penalizing candidates
 210 similar to already-selected examples, thereby promoting diversity.

211 **Complexity.** The per-candidate quantities $\{u_i\}$, $\{w_i := H_f u_i\}$, and $\{q_i := u_i^\top w_i\}$ in Algorithm 1
 212 depend on the candidate pool and target f but not on the iterates S , so we precompute them once
 213 before the greedy loop. Each iteration then reduces to $O(P)$ inner products in the dimension d
 214 of u_i plus a single update of the accumulator w , yielding a total cost of $O(KPd)$. Wall-clock
 215 measurements are reported in Appendix E.4.

216 4 Faithfulness to Retraining

217 We validate the interaction-aware estimator developed in Section 3.1 along two axes. We first measure
 218 how faithfully each estimator tracks the leave-group-out retraining effect on small-scale models,
 219 where exact ground truth can be computed, and compare against first-order influence and prior
 220 attribution methods. We then examine the pairwise interaction term’s behavior across class pairs to
 221 illustrate the class-agreement structure of Proposition 1.

222 **Setup.** We report the Spearman rank correlation between estimated and ground-truth group influ-
 223 ences. The ground-truth $\mathcal{I}^-(S)$ is obtained by retraining from scratch on $\mathcal{D} \setminus S$ and measuring the
 224 change in average held-out test loss. This evaluation is analogous to the linear datamodeling score
 225 (LDS) [43], but we avoid the term to prevent ambiguity, as our attribution is not a linear combination
 226 of individual scores.

227 We consider six dataset–model pairs spanning classification (MNIST [34], FashionMNIST [60],
 228 CIFAR-10 [32]) and regression (Concrete [63], Parkinsons [53]), with models ranging from logistic
 229 regression (LR) to MLPs and ResNet-9 [23]. To evaluate our estimator under conditions that are
 230 challenging for group attribution, we construct groups with high intra-group similarity, where existing
 231 group estimators are known to break down [25].

232 We compare our practical estimator against three influence-based methods. Throughout, we use
 233 F to denote the first-order influence term, B for the second-order term from Basu et al. [5], and
 234 I for our interaction term. The baselines are then the first-order influence function [29] (F), its
 235 second-order variant [5] ($F+B$), and our estimator augmented with the Basu term ($F+I+B$). Our
 236 method corresponds to $F+I$. We additionally compare against two attribution methods aggregated
 237 to the group level by summing individual scores: TRAK [43] and TracIn [44]. Group construction,
 238 training, and influence computation details are provided in Appendix E.

239 **Results.** Figure 1 reports the Spearman rank correlation
 240 across the six dataset–model pairs. The accuracy of first-
 241 order influence (F) varies substantially across settings, rang-
 242 ing from strong correlation on MNIST to near-zero on the
 243 harder pairs. Our estimator ($F+I$) achieves the highest correla-
 244 tion in every setting, improving upon first-order influence
 245 by up to 0.67 and yielding meaningfully positive correla-
 246 tions even where first-order influence performs poorly. This
 247 supports our claim that the interaction term captures inform-
 248 ation that first-order influence systematically overlooks.
 249 The second-order variant ($F+B$) shows instability, empiri-
 250 cally confirming the ill-conditioning of expanding $\hat{\theta}(\epsilon; S)$
 251 discussed in Section 3.1. The same correction applied to
 252 our estimator ($F+I+B$) yields either a negligible change or
 253 a substantial drop. TRAK and TracIn perform comparably
 254 to first-order influence. These findings are robust to the
 255 damping hyperparameter, as shown in Appendix F.

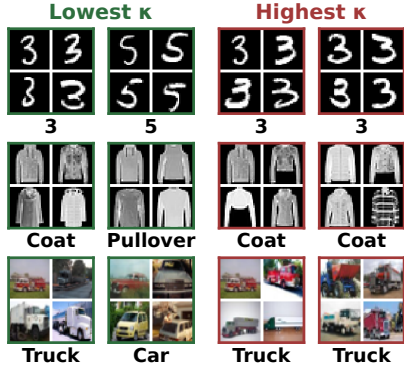


Figure 2: Representative images from the class pairs with the lowest (left) and highest (right) average pairwise interaction.

256 **Pairwise interaction analysis across classes.** Given a multiclass classification problem, we identify
 257 the class pairs whose joint effect on the test loss differs most from the sum of their individual
 258 influences. For each pair of classes (i, j) , including within-class pairs $i = j$, we compute the average
 259 pairwise interaction κ over all example pairs with one example in class i and the other in class j .
 260 Figure 2 shows representative images from the class pairs with the lowest and highest averages for
 261 MNIST/LR, FashionMNIST/MLP, and CIFAR-10/ResNet-9. The lowest- κ pairs are visually similar
 262 examples from different classes, while the highest- κ pairs come from the same class, consistent with
 263 Proposition 1.

264 5 Application to Subset Selection

265 We now evaluate our estimator as a selection criterion via Algorithm 1. Section 5.1 validates the
 266 procedure at a scale where ground-truth subset quality can be measured by retraining. Section 5.2
 267 then tests whether the gains scale to instruction-tuning data selection for Llama-3.1-8B [21].

268 5.1 Validation on Small-Scale Models

269 This subsection isolates the contribution of the interaction term to selection quality. We compare
 270 against influence-function variants on the held-out test loss after retraining, and additionally analyze
 271 the diversity of the selected subsets through their class composition.

272 **Setup.** We apply Algorithm 1 to two-layer MLPs on MNIST and FashionMNIST. For each
 273 dataset, we construct a candidate pool $\mathcal{D}_{\text{pool}}$ of $|\mathcal{D}_{\text{pool}}| = 5,000$ examples sampled uniformly
 274 at random from the training set. For each selection method, we form a subset $S \subseteq \mathcal{D}_{\text{pool}}$
 275 of size $K \in \{500, 1000, \dots, 5000\}$, retrain the MLP from scratch on S , and report the result-
 276 ing held-out test loss. To diagnose selection diversity, we additionally report the class entropy
 277 $\mathcal{H}(S) := -\sum_c p_c(S) \log p_c(S)$ of the selected subset, where $p_c(S)$ is the fraction of examples in S
 278 belonging to class c , with higher values indicating more uniform class coverage. We compare our
 279 method ($F+I$) against the same influence-function baselines as in Section 4 (F , $F+B$, $F+I+B$) and
 280 random selection. Further details are in Appendix E.

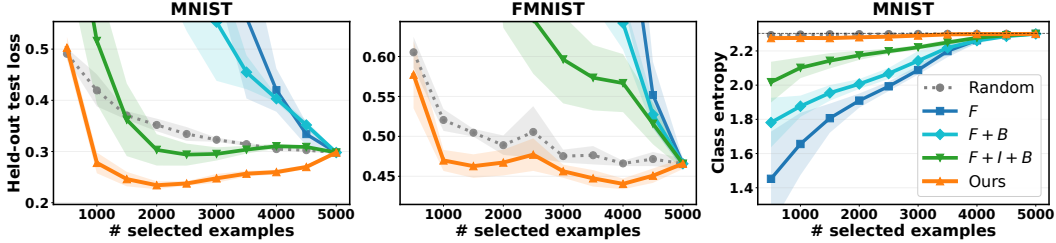


Figure 3: Held-out test loss after retraining on the selected subset on MNIST (left) and FashionMNIST (middle); class entropy of the subset on MNIST (right). Across selection sizes, lines and shaded regions show the mean and standard deviation over five seeds.

281 **Results.** The left and middle plots in Figure 3 report the held-out test loss attained by retraining on
 282 each method’s selected subset. Our method ($F+I$) outperforms every baseline on both datasets and at
 283 every selection size, indicating that the interaction term translates into substantially better subsets.
 284 The first-order influence function (F) degrades sharply and falls below random selection. We attribute
 285 this to the redundancy of the subsets it selects and analyze it in detail in the class-entropy analysis
 286 below. The second-order variant of Basu et al. [5] ($F+B$) improves over F but still trails random
 287 selection at most budgets.

288 The right plot of Figure 3 reports the class entropy of subsets selected by each method. At small
 289 budgets, both F and $F+B$ collapse onto a few classes, whereas our method matches the entropy of
 290 random selection throughout, confirming that the interaction term penalizes redundant candidates.
 291 A subset that omits entire classes cannot train a competitive classifier, which directly accounts for
 292 the test-loss degradation observed above. The corresponding plot for FashionMNIST is provided in
 293 Appendix F.

294 5.2 Data Selection for LLM Instruction Tuning

295 We now evaluate our selection method at LLM scale by fine-tuning Llama-3.1-8B on selected
 296 instruction-tuning subsets and measuring downstream task performance.

297 **Setup.** We use the LESS pool [59] of approximately 270K instruction-tuning examples as the
 298 candidate pool. Following standard practice in influence-based data selection for LLMs [59], we
 299 adopt a warmup-then-select pipeline: we first fine-tune Llama-3.1-8B on a small random subset to
 300 obtain a warmup checkpoint, which serves as $\hat{\theta}$ in our framework. We then compute per-example
 301 projected gradients u_j and target-curvature factors $H_f u_j$ at this checkpoint using LoGra [10], and
 302 run Algorithm 1 to select $K = 13,534$ examples, amounting to 5% of the pool, independently
 303 for each target task. The target function f is the instruction-masked likelihood, evaluated on the
 304 official validation set when available and on a held-out split of the training set otherwise. We then
 305 re-finetune from the original pretrained checkpoint on the selected subset. Training details, including
 306 warmup configuration, optimizer, LoRA setup, damping, and per-method runtimes, are provided in
 307 Appendix E.

308 We evaluate on seven downstream tasks covering three reasoning capabilities. For mathematical
 309 reasoning, we use GSM8K [13] and AQuA [36]. For commonsense and science reasoning, we use
 310 ARC-Easy [12], HellaSwag [66], PIQA [8], and ECQA [2]. For reading comprehension, we use
 311 SQuAD [45]. We report exact match for GSM8K, F1 for SQuAD, and accuracy for the multiple-
 312 choice tasks.

313 We compare against two influence-based baselines, two representation-similarity baselines, and a
 314 random baseline. The influence-based baselines select examples with the highest target-influence
 315 scores: Additive IF uses the standard first-order influence function [29], sharing our pipeline and
 316 differing only in the absence of the interaction term. LESS [59] uses cosine similarity between
 317 low-rank projections of training and target gradients. The representation-similarity baselines select
 318 examples whose feature embeddings most closely match those of the target task: RDS+ [27] uses
 319 pretrained-LM hidden states as features, and NV-Embed [35] uses embeddings from a retrieval-trained
 320 LLM. Random selection samples K examples uniformly from the pool.

Table 1: Test performance of fine-tuned Llama-3.1-8B on seven target tasks. We report the mean and standard deviation over five seeds; the best result in each column is shown in **bold**.

Method	GSM8K	AQuA	ARC-E	HellaSwag	PIQA	ECQA	SQuAD	Avg.
Random	46.75±2.69	24.65±0.72	90.03±0.19	56.46±0.69	77.50±0.47	69.85±0.22	75.76±0.55	63.00
Additive IF	19.35±1.07	20.63±1.03	88.72±0.30	48.66±1.99	77.38±0.79	69.52±0.29	76.92±0.30	57.31
LESS	41.83±1.32	8.82±2.87	90.11±0.15	62.11±0.87	77.40±0.49	71.19±0.12	76.73±0.51	61.17
RDS+	58.44±1.29	24.65±2.87	88.70±0.32	52.19±1.03	75.14±1.14	70.88±0.41	57.60±1.91	61.09
NV-Embed	57.97±0.47	8.43±1.06	89.01±0.16	53.94±0.74	75.13±0.44	71.75±0.24	76.52±1.10	61.82
Ours	52.69±0.96	30.94±0.72	90.44±0.11	63.43±0.34	79.39±0.32	71.17±0.19	79.87±0.34	66.85

321 **Results.** Table 1 reports the mean and standard deviation over five seeds. Our method attains
 322 the best performance on five of the seven tasks and surpasses random selection on every task. The
 323 two influence-based baselines, Additive IF and LESS, are outperformed on six of the seven tasks;
 324 Additive IF, which ablates our interaction term, even falls below random on most tasks, isolating
 325 the contribution of interaction-aware selection. LESS and NV-Embed further drop to below 9%
 326 accuracy on AQuA, whereas our method remains stable across all tasks. On GSM8K, although our
 327 method outperforms random selection and the influence-based baselines, the representation-similarity
 328 baselines perform especially well. We conjecture this stems from GSM8K’s narrow design: problems
 329 are restricted to elementary arithmetic, single integer answers, and short, simple descriptions. This
 330 concentrated distribution plays to the strength of representation-similarity selection, which directly
 331 retrieves semantically near-identical examples from the pool.

332 6 Related Work

333 **Data attribution and influence functions.** Data attribution methods quantify the contribution of
 334 individual training examples to a model’s learned parameters or predictions on a target [20, 29, 43, 44].
 335 Among these, influence functions [29] approximate this counterfactual via local linearization, with
 336 subsequent work improving scalability through random-projection approaches [43], trajectory-based
 337 methods [44], efficient curvature approximations [22, 49, 54], and extensions to large language
 338 and generative models [1, 9, 10, 33, 40]. Complementary work studies the fragility of influence
 339 estimates [6, 17, 31, 47]. Ye et al. [62] also expand the validation loss to second order, but only for
 340 single-example influence in noisy-label settings, without group attribution. Prior work also shows
 341 that summing individual attributions mismeasures group-level effects [25, 30, 48]. Closest to our
 342 setting, Basu et al. [5] address this additive gap via a response-side expansion; we instead expand the
 343 target function, yielding closed-form pairwise interactions. Wang et al. [56] and Wei et al. [58] also
 344 model joint influence but focus on trajectory-accumulated rather than counterfactual effects.

345 **Training data selection.** Data selection methods score examples by influence [1, 10, 16, 59], diver-
 346 sity [42], learned utility [18, 64], or Shapley-style contributions [20]. Scalable variants leverage influ-
 347 ence distillation or small-model trajectories [41, 61]. Redundancy-aware pruning further penalizes
 348 similar examples via pairwise similarity objectives [51]. However, individually influential examples
 349 need not form the best subset once interactions make collective effects non-additive [25, 26, 55, 65].
 350 Our second-order expansion of the target function enables interaction-aware selection by updating
 351 each candidate’s marginal utility against the already-selected set.

352 7 Conclusion

353 We presented an interaction-aware influence function obtained from a second-order expansion of
 354 the target function, decomposing group attribution into a standard first-order term and a pairwise
 355 interaction term. Empirically, our method improves Spearman correlation with ground-truth retraining
 356 effects on six dataset–model pairs. In subset selection, it outperforms random selection on every task
 357 and beats existing influence-based and representation-similarity baselines on five of seven downstream
 358 tasks. Two limitations suggest natural directions for future work: our greedy selection performs
 359 strongly in practice but lacks formal optimality guarantees, and our estimator inherits approximation
 360 error from the underlying Hessian approximations. Advances in scalable inverse-Hessian estimation
 361 would directly improve robustness.

362 References

- 363 [1] Ishika Agarwal and Dilek Hakkani-Tür. Neural networks for learnable and scalable influence
364 estimation of instruction fine-tuning data. In *The Thirty-ninth Annual Conference on Neural*
365 *Information Processing Systems, 2025*.
- 366 [2] Shourya Aggarwal, Divyanshu Mandowara, Vishwajeet Agrawal, Dinesh Khandelwal, Parag
367 Singla, and Dinesh Garg. Explanations for commonsenseqa: New dataset and models. In
368 *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and*
369 *the 11th International Joint Conference on Natural Language Processing, 2021*.
- 370 [3] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agar-
371 wal. Deep batch active learning by diverse, uncertain gradient lower bounds. *International*
372 *Conference on Learning Representations, 2020*.
- 373 [4] Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger B Grosse. If influence
374 functions are the answer, then what is the question? *Advances in Neural Information Processing*
375 *Systems, 2022*.
- 376 [5] Samyadeep Basu, Xuchen You, and Soheil Feizi. On second-order group influence functions
377 for black-box predictions. In *International Conference on Machine Learning, 2020*.
- 378 [6] Samyadeep Basu, Philip Pope, and Soheil Feizi. Influence functions in deep learning are fragile.
379 *International Conference on Learning Representations, 2021*.
- 380 [7] Vighnesh Birodkar, Hossein Mobahi, and Samy Bengio. Semantic redundancies in image-
381 classification datasets: The 10% you don't need. *arXiv preprint arXiv:1901.11409, 2019*.
- 382 [8] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning
383 about physical commonsense in natural language. In *Proceedings of the AAAI Conference on*
384 *Artificial Intelligence, 2020*.
- 385 [9] Tyler A. Chang, Dheeraj Rajagopal, Tolga Bolukbasi, Lucas Dixon, and Ian Tenney. Scalable
386 influence and fact tracing for large language model pretraining. In *The Thirteenth International*
387 *Conference on Learning Representations, 2025*.
- 388 [10] Sang Keun Choe, Hwijeen Ahn, Juhan Bae, Kewen Zhao, Minsoo Kang, Youngseog Chung,
389 Adithya Pratapa, Willie Neiswanger, Emma Strubell, Teruko Mitamura, et al. What is your data
390 worth to gpt? llm-scale data valuation with influence functions. *Advances in neural information*
391 *processing systems, 2025*.
- 392 [11] Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin
393 Rostamizadeh, and Sanjiv Kumar. Batch active learning at scale. *Advances in Neural Information*
394 *Processing Systems, 2021*.
- 395 [12] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick,
396 and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning
397 challenge. *arXiv preprint arXiv:1803.05457, 2018*.
- 398 [13] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
399 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
400 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168,*
401 *2021*.
- 402 [14] Cody Coleman, Deepak Narayanan, Daniel Kang, Tian Zhao, Jian Zhang, Luigi Nardi, Peter
403 Bailis, Kunle Olukotun, Chris Ré, and Matei Zaharia. Dawnbench: An end-to-end deep learning
404 benchmark and competition. *Training, 2017*.
- 405 [15] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning, 1995*.
- 406 [16] Qirun Dai, Dylan Zhang, Jiaqi W Ma, and Hao Peng. Improving influence-based instruction
407 tuning data selection for balanced learning of diverse capabilities. *Findings of the Association*
408 *for Computational Linguistics, 2025*.

- 409 [17] Junwei Deng, Weijing Tang, and Jiaqi W. Ma. A versatile influence function for data attribution
410 with non-decomposable loss. *arXiv preprint arXiv:2412.01335*, 2024.
- 411 [18] Logan Engstrom, Axel Feldmann, and Aleksander Madry. Dsdm: Model-aware dataset selection
412 with datamodels. In *International Conference on Machine Learning*, 2024.
- 413 [19] Thomas George, César Laurent, Xavier Bouthillier, Nicolas Ballas, and Pascal Vincent. Fast
414 approximate natural gradient descent in a kronecker factored eigenbasis. *Advances in neural
415 information processing systems*, 2018.
- 416 [20] Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine
417 learning. In *International conference on machine learning*, 2019.
- 418 [21] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian,
419 Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama
420 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 421 [22] Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini,
422 Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al. Studying large language model
423 generalization with influence functions. *arXiv preprint arXiv:2308.03296*, 2023.
- 424 [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
425 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
426 2016.
- 427 [24] Jaeseung Heo, Kyeongheung Yun, Seokwon Yoon, MoonJeong Park, Jungseul Ok, and Dong-
428 woo Kim. Influence functions for edge edits in non-convex graph neural networks. *Advances in
429 Neural Information Processing Systems*, 2025.
- 430 [25] Yuzheng Hu, Pingbang Hu, Han Zhao, et al. Most influential subset selection: Challenges,
431 promises, and beyond. *Advances in Neural Information Processing Systems*, 2024.
- 432 [26] Jenny Y Huang, David R Burt, Yunyi Shen, Tin D Nguyen, and Tamara Broderick. Approx-
433 imations to worst-case data dropping: unmasking failure modes. *Transactions on Machine
434 Learning Research*, 2025.
- 435 [27] Hamish Ivison, Muru Zhang, Faeze Brahman, Pang Wei Koh, and Pradeep Dasigi. Large-Scale
436 Data Selection for Instruction Tuning. *arXiv preprint arXiv:2503.01807*, 2025.
- 437 [28] Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch
438 acquisition for deep bayesian active learning. *Advances in neural information processing
439 systems*, 2019.
- 440 [29] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions.
441 In *International conference on machine learning*, 2017.
- 442 [30] Pang Wei W Koh, Kai-Siang Ang, Hubert Teo, and Percy S Liang. On the accuracy of influence
443 functions for measuring group effects. *Advances in neural information processing systems*,
444 2019.
- 445 [31] Philipp Alexander Kreer, Wilson Wu, Maxwell Adam, Zach Furman, and Jesse Hoogland.
446 Bayesian influence functions for hessian-free data attribution. In *The Fourteenth International
447 Conference on Learning Representations*, 2026.
- 448 [32] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
449 2009.
- 450 [33] Yongchan Kwon, Eric Wu, Kevin Wu, and James Zou. Datainf: Efficiently estimating data
451 influence in lora-tuned llms and diffusion models. *International Conference on Learning
452 Representations*, 2024.
- 453 [34] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning
454 applied to document recognition. *Proceedings of the IEEE*, 2002.

- 455 [35] Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan
456 Catanzaro, and Wei Ping. Nv-embed: Improved techniques for training llms as generalist
457 embedding models. *International Conference on Learning Representations*, 2025.
- 458 [36] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale
459 generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th*
460 *Annual Meeting of the Association for Computational Linguistics*, 2017.
- 461 [37] James Martens. New insights and perspectives on the natural gradient method. *Journal of*
462 *Machine Learning Research*, 2020.
- 463 [38] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approx-
464 imate curvature. In *International conference on machine learning*, 2015.
- 465 [39] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of
466 machine learning models. In *International Conference on Machine Learning*, 2020.
- 467 [40] Bruno Kacper Mlodozieniec, Runa Eschenhagen, Juhan Bae, Alexander Immer, David Krueger,
468 and Richard E. Turner. Influence functions for scalable data attribution in diffusion models. In
469 *The Thirteenth International Conference on Learning Representations*, 2025.
- 470 [41] Mahdi Nikdan, Vincent Cohen-Addad, Dan Alistarh, and Vahab Mirrokni. Efficient data
471 selection at scale via influence distillation. In *The Thirty-ninth Annual Conference on Neural*
472 *Information Processing Systems*, 2025.
- 473 [42] Xingyuan Pan, Luyang Huang, Liyan Kang, Zhicheng Liu, Yu Lu, and Shanbo Cheng. G-
474 dig: Towards gradient-based diverse and high-quality instruction data selection for machine
475 translation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational*
476 *Linguistics (Volume 1: Long Papers)*, 2024.
- 477 [43] Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry.
478 Trak: Attributing model behavior at scale. In *International Conference on Machine Learning*,
479 2023.
- 480 [44] Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training
481 data influence by tracing gradient descent. *Advances in Neural Information Processing Systems*,
482 2020.
- 483 [45] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions
484 for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical*
485 *Methods in Natural Language Processing*, 2016.
- 486 [46] Joshua David Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. Contrastive
487 learning with hard negative samples. In *International Conference on Learning Representations*,
488 2021.
- 489 [47] Ittai Rubinstein and Samuel B. Hopkins. Rescaled influence functions: Accurate data attribution
490 in high dimension. In *The Thirty-ninth Annual Conference on Neural Information Processing*
491 *Systems*, 2025.
- 492 [48] Nikunj Saunshi, Arushi Gupta, Mark Braverman, and Sanjeev Arora. Understanding influ-
493 ence functions and datamodels via harmonic analysis. *International Conference on Learning*
494 *Representations*, 2023.
- 495 [49] Andrea Schioppa, Polina Zablotskaia, David Vilar, and Artem Sokolov. Scaling up influence
496 functions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
- 497 [50] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors
498 with online hard example mining. In *Proceedings of the IEEE conference on computer vision*
499 *and pattern recognition*, 2016.
- 500 [51] Haoru Tan, Sitong Wu, Wei Huang, Shizhen Zhao, and XIAOJUAN QI. Data pruning by infor-
501 mation maximization. In *The Thirteenth International Conference on Learning Representations*,
502 2025.

- 503 [52] Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio,
504 and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network
505 learning. In *International Conference on Learning Representations*, 2019.
- 506 [53] Athanasios Tsanas, Max Little, Patrick McSharry, and Lorraine Ramig. Accurate telemonitoring
507 of parkinson’s disease progression by non-invasive speech tests. *Nature Precedings*, 2009.
- 508 [54] Andrew Wang, Elisa Nguyen, Runshi Yang, Juhan Bae, Sheila A McIlraith, and Roger Grosse.
509 Better training data attribution via better inverse hessian-vector products. *Advances in Neural
510 Information Processing Systems*, 2025.
- 511 [55] Jiachen T Wang, Tianji Yang, James Zou, Yongchan Kwon, and Ruoxi Jia. Rethinking data
512 shapley for data selection tasks: Misleads and merits. *International Conference on Machine
513 Learning*, 2024.
- 514 [56] Jiachen T Wang, Prateek Mittal, Dawn Song, and Ruoxi Jia. Data shapley in one training run.
515 *International Conference on Learning Representations*, 2025.
- 516 [57] Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu,
517 David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go?
518 exploring the state of instruction tuning on open resources. *Advances in Neural Information
519 Processing Systems*, 2023.
- 520 [58] Jingyu Wei, Bo Liu, Tianjiao Wan, Baoyun Peng, Xingkong Ma, and Mengmeng Guo. Ji2s:
521 Joint influence-aware instruction data selection for efficient fine-tuning. In *Proceedings of the
522 2025 Conference on Empirical Methods in Natural Language Processing*, 2025.
- 523 [59] Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. LESS:
524 Selecting influential data for targeted instruction tuning. In *International Conference on Machine
525 Learning*, 2024.
- 526 [60] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
527 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 528 [61] Yu Yang, Siddhartha Mishra, Jeffrey N Chiang, and Baharan Mirzasoleiman. Smalltolarge
529 (s2l): Scalable data selection for fine-tuning large language models by summarizing training
530 trajectories of small models. In *The Thirty-eighth Annual Conference on Neural Information
531 Processing Systems*, 2024.
- 532 [62] Xichen Ye, Yifan Wu, Weizhong Zhang, Cheng Jin, and Yifan Chen. Towards robust influence
533 functions with flat validation minima. *arXiv preprint arXiv:2505.19097*, 2025.
- 534 [63] I-C Yeh. Modeling of strength of high-performance concrete using artificial neural networks.
535 *Cement and Concrete research*, 1998.
- 536 [64] Zichun Yu, Spandan Das, and Chenyan Xiong. Mates: Model-aware data selection for efficient
537 pretraining with data influence models. *Advances in Neural Information Processing Systems*,
538 2024.
- 539 [65] Zichun Yu, Fei Peng, Jie Lei, Arnold Overwijk, Wen tau Yih, and Chenyan Xiong. Group-level
540 data selection for efficient pretraining. In *The Thirty-ninth Annual Conference on Neural
541 Information Processing Systems*, 2025.
- 542 [66] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can
543 a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the
544 Association for Computational Linguistics*, 2019.
- 545 [67] Chi Zhang, Huaping Zhong, Kuan Zhang, Chengliang Chai, Rui Wang, Xinlin Zhuang, Tianyi
546 Bai, Jiantao Qiu, Lei Cao, Ju Fan, et al. Harnessing diversity for important data selection in
547 pretraining large language models. *International Conference on Learning Representations*,
548 2025.

550 A Notation

551 Table 2 consolidates the notation used throughout the paper. The symbols are grouped into five
 552 categories: data and model objects, group-level effects on the target function, parameter shifts and
 553 curvature matrices, the pairwise interaction term, and quantities specific to data selection. Within
 554 each category, definitions are listed in the order they first appear in the main text, and we use the
 555 same symbols in the appendix derivations.

Table 2: Summary of notation used throughout the main text and appendices, organized by category.

Symbol	Description
Data and model	
$\mathcal{D} = \{z_i\}_{i=1}^N$	Training dataset of N examples
$z_i = (x_i, y_i)$	i -th training example with input x_i and label y_i
$S \subseteq \mathcal{D}$	Subset (group) of training examples
$\theta \in \mathbb{R}^p$	Model parameters
$\hat{\theta}$	Parameters obtained by training on \mathcal{D}
$\hat{\theta}_{\mathcal{D} \setminus S}$	Parameters obtained by retraining on $\mathcal{D} \setminus S$
$\hat{\theta}_{\mathcal{D} \cup S}$	Parameters obtained by retraining on $\mathcal{D} \cup S$
$\hat{\theta}(\epsilon; S)$	Parameters when examples in S are upweighted by ϵ
$\ell(z, \theta)$	Per-example loss
$L(\theta) = \frac{1}{N} \sum_i \ell(z_i, \theta)$	Empirical training loss
$f(\theta)$	Target function (e.g., loss on a held-out example)
Group-level effects	
$\mathcal{I}^-(S)$	Removal effect: $f(\hat{\theta}_{\mathcal{D} \setminus S}) - f(\hat{\theta})$
$\mathcal{I}^+(S)$	Inclusion effect: $f(\hat{\theta}_{\mathcal{D} \cup S}) - f(\hat{\theta})$
$\hat{\mathcal{I}}_{\text{lin}}^-(S)$	First-order estimate of $\mathcal{I}^-(S)$
$\hat{\mathcal{I}}_{\text{lin}}^+(S)$	First-order estimate of $\mathcal{I}^+(S)$
$\hat{\mathcal{I}}^-(S)$	Our interaction-aware estimator of $\mathcal{I}^-(S)$
$\hat{\mathcal{I}}^+(S)$	Our interaction-aware estimator of $\mathcal{I}^+(S)$
Parameter shifts and curvature	
$g_i := \nabla_{\theta} \ell(z_i, \hat{\theta})$	Per-example gradient at $\hat{\theta}$
$H := \nabla_{\theta}^2 L(\hat{\theta})$	Training-loss Hessian at $\hat{\theta}$
G	Gauss–Newton Hessian of the training loss at $\hat{\theta}$
$H_f := \nabla_{\theta}^2 f(\hat{\theta})$	Target-function Hessian at $\hat{\theta}$
$M := H^{-1} H_f H^{-1}$	Bilinear form arising in Proposition 1
$u_i := H^{-1} g_i$	Per-example parameter shift
$u_S := \sum_{z_i \in S} u_i$	Aggregate parameter shift over S
Pairwise interaction term	
$\kappa(z_i, z_j) := u_i^{\top} H_f u_j$	Pairwise interaction between z_i and z_j
Data selection	
$\mathcal{D}_{\text{pool}}$	Candidate pool of examples
K	Selection budget (target subset size)
S_t	Selected subset after t greedy iterations
$m(z_j S)$	Marginal score of adding z_j to current subset S

556 B Derivation of the First-Order Influence Function

557 We derive Equation (4) and Equation (5), the first-order Taylor approximations of (i) the parameter
 558 shift induced by removing S and (ii) the resulting change in the target function f . Throughout, we
 559 maintain the assumptions stated in Section 2: $\ell(z_i, \cdot)$ is twice continuously differentiable for every

560 $z_i \in \mathcal{D}$, and the training-loss Hessian

$$H := \frac{1}{N} \sum_{z_i \in \mathcal{D}} \nabla_{\theta}^2 \ell(z_i, \hat{\theta})$$

561 is nonsingular at $\hat{\theta}$, and the target function f is differentiable at $\hat{\theta}$.

562 **Proof strategy.** Our goal is to approximate the target change $f(\hat{\theta}_{\mathcal{D} \setminus S}) - f(\hat{\theta})$ caused by removing S ,
 563 which we obtain in two first-order Taylor-expansion steps. The first step approximates the parameter
 564 shift $\hat{\theta}_{\mathcal{D} \setminus S} - \hat{\theta}$: as noted in Section 2, this parameter shift corresponds to a finite perturbation from
 565 $\epsilon = 0$ to $\epsilon = -1/N$ along the reweighting path $\hat{\theta}(\epsilon; S)$, and we approximate it by a first-order Taylor
 566 expansion of $\hat{\theta}(\epsilon; S)$ in ϵ , yielding Equation (4). The second step substitutes the resulting parameter
 567 shift into a first-order Taylor expansion of f around $\hat{\theta}$, yielding Equation (5). The only quantity we
 568 need to compute for the first step is the slope $d\hat{\theta}(\epsilon; S)/d\epsilon$ at $\epsilon = 0$. To compute it, we use the fact
 569 that $\hat{\theta}(\epsilon; S)$ minimizes the perturbed objective for every ϵ , so the gradient of the perturbed objective
 570 at $\hat{\theta}(\epsilon; S)$ remains zero as ϵ varies. Differentiating this zero-gradient identity in ϵ yields a closed form
 571 for the slope.

572 **Setup.** Recall the perturbed objective from Equation (3),

$$\mathcal{R}(\theta; \epsilon, S) := \frac{1}{N} \sum_{z_i \in \mathcal{D}} \ell(z_i, \theta) + \epsilon \sum_{z_i \in S} \ell(z_i, \theta), \quad (13)$$

573 and its minimizer $\hat{\theta}(\epsilon; S) := \arg \min_{\theta} \mathcal{R}(\theta; \epsilon, S)$. At $\epsilon = 0$ the perturbation vanishes and we recover
 574 $\hat{\theta}(0; S) = \hat{\theta}$. At $\epsilon = -1/N$, substituting into Equation (13) cancels the loss terms for $z_i \in S$:

$$\mathcal{R}(\theta; -\frac{1}{N}, S) = \frac{1}{N} \sum_{z_i \in \mathcal{D}} \ell(z_i, \theta) - \frac{1}{N} \sum_{z_i \in S} \ell(z_i, \theta) = \frac{1}{N} \sum_{z_i \in \mathcal{D} \setminus S} \ell(z_i, \theta).$$

575 This differs from the empirical risk $\mathcal{R}_{\mathcal{D} \setminus S}(\theta) = \frac{1}{|\mathcal{D} \setminus S|} \sum_{z_i \in \mathcal{D} \setminus S} \ell(z_i, \theta)$ defined in Section 2 only by
 576 a positive constant factor, so its minimizer coincides with the leave-group-out parameter $\hat{\theta}_{\mathcal{D} \setminus S}$.

577 **Zero-gradient identity at the minimizer.** Since $\hat{\theta}(\epsilon; S)$ is a minimizer of $\mathcal{R}(\cdot; \epsilon, S)$, the gradient
 578 of \mathcal{R} in θ is zero at $\hat{\theta}(\epsilon; S)$. Concretely, computing $\nabla_{\theta} \mathcal{R}(\theta; \epsilon, S)$ from Equation (13) and evaluating
 579 at $\theta = \hat{\theta}(\epsilon; S)$ gives

$$\frac{1}{N} \sum_{z_i \in \mathcal{D}} \nabla_{\theta} \ell(z_i, \hat{\theta}(\epsilon; S)) + \epsilon \sum_{z_i \in S} \nabla_{\theta} \ell(z_i, \hat{\theta}(\epsilon; S)) = 0. \quad (14)$$

580 This identity holds for every ϵ in a neighborhood of 0. By the implicit function theorem, the
 581 nonsingularity of H further ensures that the minimizer path $\hat{\theta}(\epsilon; S)$ is differentiable in ϵ near $\epsilon = 0$,
 582 so we may differentiate both sides of the identity in ϵ to extract the desired slope $d\hat{\theta}(\epsilon; S)/d\epsilon|_{\epsilon=0}$.

583 **Differentiating in ϵ .** Each per-example gradient $\nabla_{\theta} \ell(z_i, \hat{\theta}(\epsilon; S))$ depends on ϵ only through $\hat{\theta}(\epsilon; S)$.
 584 Differentiating Equation (14) in ϵ and applying the chain rule to each such term gives

$$\underbrace{\left[\frac{1}{N} \sum_{z_i \in \mathcal{D}} \nabla_{\theta}^2 \ell(z_i, \hat{\theta}(\epsilon; S)) + \epsilon \sum_{z_i \in S} \nabla_{\theta}^2 \ell(z_i, \hat{\theta}(\epsilon; S)) \right]}_{\text{Hessian of } \mathcal{R}(\cdot; \epsilon, S)} \frac{d\hat{\theta}(\epsilon; S)}{d\epsilon} + \sum_{z_i \in S} \nabla_{\theta} \ell(z_i, \hat{\theta}(\epsilon; S)) = 0. \quad (15)$$

585 At $\epsilon = 0$, the second summand inside the bracket is zero due to its ϵ prefactor, and the bracket reduces
 586 to H . Using $g_i := \nabla_{\theta} \ell(z_i, \hat{\theta})$ as defined in Section 2, Equation (15) becomes

$$H \cdot \frac{d\hat{\theta}(\epsilon; S)}{d\epsilon} \Big|_{\epsilon=0} + \sum_{z_i \in S} g_i = 0. \quad (16)$$

587 Since H is invertible, we solve for the slope:

$$\frac{d\hat{\theta}(\epsilon; S)}{d\epsilon} \Big|_{\epsilon=0} = -H^{-1} \sum_{z_i \in S} g_i. \quad (17)$$

588 **Parameter-shift approximation (Equation (4)).** We now use the slope in Equation (17) inside the
 589 first Taylor expansion promised in the strategy. Expanding $\hat{\theta}(\epsilon; S)$ to first order around $\epsilon = 0$ gives

$$\hat{\theta}(\epsilon; S) \approx \hat{\theta} + \epsilon \cdot \left. \frac{d\hat{\theta}(\epsilon; S)}{d\epsilon} \right|_{\epsilon=0} = \hat{\theta} - \epsilon H^{-1} \sum_{z_i \in S} g_i. \quad (18)$$

590 Setting $\epsilon = -1/N$, which corresponds to leaving S out, recovers Equation (4):

$$\hat{\theta}_{\mathcal{D} \setminus S} - \hat{\theta} \approx \frac{1}{N} H^{-1} \sum_{z_i \in S} g_i. \quad (19)$$

591 **Target-change approximation via the parameter shift (Equation (5)).** For the second step, a
 592 first-order Taylor expansion of the target f around $\hat{\theta}$ gives

$$f(\hat{\theta}_{\mathcal{D} \setminus S}) - f(\hat{\theta}) \approx \nabla_{\theta} f(\hat{\theta})^{\top} (\hat{\theta}_{\mathcal{D} \setminus S} - \hat{\theta}). \quad (20)$$

593 Substituting the parameter shift from Equation (19) yields the standard first-order group influence
 594 estimate of Equation (5):

$$\hat{\mathcal{I}}_{\text{lin}}^{-}(S) := \frac{1}{N} \nabla_{\theta} f(\hat{\theta})^{\top} H^{-1} \sum_{z_i \in S} g_i. \quad (21)$$

595 The additivity follows immediately, since the right-hand side depends on S only through the linear
 596 sum $\sum_{z_i \in S} g_i$.

597 C Derivations and Proofs for the Interaction-Aware Influence Function

598 C.1 Derivation of the interaction-aware influence function

599 We derive the interaction-aware influence function in Equation (7) by combining the second-order
 600 Taylor expansion of the target function with the first-order parameter shift induced by removing S .

601 Recall from Equation (4) in Section 2 that the parameter shift induced by removing S admits the
 602 first-order approximation

$$\delta_S := \hat{\theta}_{\mathcal{D} \setminus S} - \hat{\theta} \approx \frac{1}{N} H^{-1} \sum_{z_i \in S} g_i = \frac{1}{N} u_S, \quad (22)$$

603 where the last equality uses $u_i := H^{-1} g_i$ and $u_S := \sum_{z_i \in S} u_i$. The second-order Taylor expansion
 604 of f around $\hat{\theta}$ in Equation (6) is

$$f(\hat{\theta}_{\mathcal{D} \setminus S}) - f(\hat{\theta}) = \nabla_{\theta} f(\hat{\theta})^{\top} \delta_S + \frac{1}{2} \delta_S^{\top} H_f \delta_S + O(\|\delta_S\|^3). \quad (23)$$

605 Substituting Equation (22) into the linear term of Equation (23) gives

$$\nabla_{\theta} f(\hat{\theta})^{\top} \delta_S \approx \frac{1}{N} \nabla_{\theta} f(\hat{\theta})^{\top} u_S. \quad (24)$$

606 Substituting the same approximation into the quadratic term gives

$$\frac{1}{2} \delta_S^{\top} H_f \delta_S \approx \frac{1}{2} \left(\frac{1}{N} u_S \right)^{\top} H_f \left(\frac{1}{N} u_S \right) = \frac{1}{2N^2} u_S^{\top} H_f u_S, \quad (25)$$

607 where the factor $1/N^2$ arises from the two factors of $1/N$ in δ_S . Combining Equations (24) and (25)
 608 and recalling that the leave-group-out effect is $\mathcal{I}^{-}(S) = f(\hat{\theta}_{\mathcal{D} \setminus S}) - f(\hat{\theta})$ yields

$$\mathcal{I}^{-}(S) \approx \frac{1}{N} \nabla_{\theta} f(\hat{\theta})^{\top} u_S + \frac{1}{2N^2} u_S^{\top} H_f u_S, \quad (26)$$

609 which is Equation (7) in the main text. The remainder term $O(\|\delta_S\|^3)$ is dropped, consistent with
 610 the second-order approximation. The first term recovers the standard first-order estimate $\hat{\mathcal{I}}_{\text{lin}}^{-}(S)$
 611 from Equation (5), and the second term is the curvature-induced interaction term that gives rise to
 612 non-additivity across examples in S .

613 **C.2 Derivation for the data-addition setting**

614 We now derive the addition-setting influence function in Equation (8). The argument parallels
 615 Appendix C.1, with the only change appearing in the sign of the parameter shift induced by adding S
 616 rather than removing it.

617 Consider the reweighted objective in Equation (3) with weight ϵ assigned to the examples in S .
 618 Setting $\epsilon = -1/N$ removes the contribution of S , recovering the leave-group-out parameter $\hat{\theta}_{\mathcal{D}\setminus S}$
 619 as a finite step along the reweighting path; setting $\epsilon = +1/N$ instead adds an extra copy of each
 620 example in S with weight $1/N$, which corresponds to retraining on the augmented set $\mathcal{D} \cup S$ with
 621 each new example carrying the same per-example weight as in \mathcal{D} . The first-order parameter shift
 622 induced by this addition is therefore

$$\hat{\theta}_{\mathcal{D}\cup S} - \hat{\theta} \approx -\frac{1}{N} H^{-1} \sum_{z_i \in S} g_i = -\frac{1}{N} u_S, \quad (27)$$

623 which differs from Equation (22) only in sign. A complete derivation through the implicit function
 624 theorem is given in Appendix B; the sign flip reflects the opposite direction of the reweighting step.

625 Substituting Equation (27) into the second-order Taylor expansion of f around $\hat{\theta}$ yields

$$\begin{aligned} f(\hat{\theta}_{\mathcal{D}\cup S}) - f(\hat{\theta}) &\approx \nabla_{\theta} f(\hat{\theta})^{\top} \left(-\frac{1}{N} u_S \right) + \frac{1}{2} \left(-\frac{1}{N} u_S \right)^{\top} H_f \left(-\frac{1}{N} u_S \right) \\ &= -\frac{1}{N} \nabla_{\theta} f(\hat{\theta})^{\top} u_S + \frac{1}{2N^2} u_S^{\top} H_f u_S. \end{aligned} \quad (28)$$

626 The linear term flips sign because the parameter shift itself is negated, while the quadratic term is
 627 invariant under this negation: the two factors of $-1/N$ in the bilinear form combine to a positive
 628 $1/N^2$. Recalling that $\mathcal{I}^+(S) = f(\hat{\theta}_{\mathcal{D}\cup S}) - f(\hat{\theta})$ gives Equation (8) in the main text.

629 This sign asymmetry has a direct consequence for our selection criterion in Section 3.2. When f is
 630 to be minimized, a more negative $\mathcal{I}^+(S)$ is preferable, so the linear term rewards candidates whose
 631 parameter shifts align with $-\nabla_{\theta} f(\hat{\theta})$. The quadratic term, in contrast, contributes the same sign in
 632 both removal and addition settings, reflecting the fact that interactions among examples in S depend
 633 on their joint geometry rather than on the direction of the perturbation.

634 **C.3 Pairwise decomposition of the interaction term**

635 The interaction term $u_S^{\top} H_f u_S$ in Equations (7) and (8) admits a pairwise decomposition used in
 636 Section 3.1, which we derive below.

637 **Pairwise decomposition.** Recall that $u_S = \sum_{z_i \in S} u_i$ is the aggregate of per-example shifts.
 638 Substituting this definition into the bilinear form and expanding gives

$$u_S^{\top} H_f u_S = \left(\sum_{z_i \in S} u_i \right)^{\top} H_f \left(\sum_{z_j \in S} u_j \right) = \sum_{z_i \in S} \sum_{z_j \in S} u_i^{\top} H_f u_j = \sum_{z_i, z_j \in S} \kappa(z_i, z_j), \quad (29)$$

639 where $\kappa(a, b) := u_a^{\top} H_f u_b$ as defined in Equation (9). This recovers the pairwise form used
 640 throughout Section 3.1.

641 **Self and cross contributions.** The double sum in Equation (29) ranges over ordered pairs and can
 642 be further separated into diagonal and off-diagonal contributions:

$$u_S^{\top} H_f u_S = \underbrace{\sum_{z_i \in S} \kappa(z_i, z_i)}_{\text{self contribution}} + 2 \underbrace{\sum_{\{z_i, z_j\} \subset S, i \neq j} \kappa(z_i, z_j)}_{\text{cross contribution}}, \quad (30)$$

643 where the cross contribution sums over unordered pairs $\{z_i, z_j\}$ with $i \neq j$ and the factor of two
 644 reflects the symmetry $\kappa(z_i, z_j) = \kappa(z_j, z_i)$ inherited from $H_f = H_f^{\top}$. The self contribution
 645 $\sum_i \kappa(z_i, z_i) = \sum_i u_i^{\top} H_f u_i$ is determined entirely by the individual examples and would persist

646 even if examples in S were processed independently. The cross contribution is the part of the
 647 interaction term that genuinely encodes interactions between distinct examples and is the source of
 648 non-additivity beyond the additive baseline. In the main text we present the unified double-sum form
 649 in Equation (9) so that the analysis remains agnostic to whether self or cross terms dominate. The
 650 split into self and cross contributions in Equation (30) becomes useful when isolating the genuinely
 651 interactional content, which we exploit in the marginal score derivation in Appendix C.7.

652 C.4 Spectral interpretation of the pairwise interaction

653 We provide a spectral interpretation of the pairwise interaction $\kappa(a, b)$ that holds in the general
 654 setting where the target Hessian H_f may be indefinite, going beyond the inner-product reading of
 655 Section 3.1.

656 **Spectral decomposition of the interaction term.** Since f is twice continuously differentiable at $\hat{\theta}$,
 657 H_f is symmetric and admits the eigendecomposition

$$H_f = \sum_k \mu_k v_k v_k^\top, \quad (31)$$

658 where $\{\mu_k\}$ are the real eigenvalues and $\{v_k\}$ form an orthonormal basis of unit eigenvectors.
 659 Substituting Equation (31) into the definition of the interaction term $\kappa(a, b) = u_a^\top H_f u_b$ yields

$$\kappa(a, b) = \sum_k \mu_k (v_k^\top u_a)(v_k^\top u_b). \quad (32)$$

660 Each summand isolates the contribution of a single eigendirection v_k . The product $(v_k^\top u_a)(v_k^\top u_b)$
 661 measures how the parameter shifts of a and b are aligned along v_k . The eigenvalue μ_k then specifies
 662 how strongly this alignment contributes to $\kappa(a, b)$, as well as its sign. The spectral view thus provides
 663 additional information beyond the bilinear-form definition, allowing $\kappa(a, b)$ to be read as a sum of
 664 independent contributions, one per eigendirection of H_f .

665 **Reading the decomposition as a weighted similarity score.** A useful way to read Equation (32) is
 666 as an aggregate similarity score. The shifts u_a and u_b are compared from multiple perspectives, one
 667 per eigendirection v_k . These per-direction comparisons are then combined through a weighted sum,
 668 with the eigenvalues μ_k acting as importance weights. When H_f is positive definite all weights μ_k
 669 are positive, so alignment along any eigendirection contributes consistently in the positive direction
 670 to the aggregate score. This recovers the inner-product reading of Section 3.1 in spectral form. The
 671 additional insight is that the weight assigned to each perspective is precisely the curvature of f along
 672 that direction.

673 **Extension to the indefinite case.** When H_f is instead indefinite, $\kappa(a, b)$ is no longer an inner
 674 product but a symmetric bilinear form. The eigendecomposition in Equation (31) remains valid, so
 675 the per-direction decomposition of $\kappa(a, b)$ in Equation (32) still applies. The only change is that
 676 some eigenvalues μ_k may now be negative. Concretely, suppose u_a and u_b are aligned along an
 677 eigendirection v_k , that is, $(v_k^\top u_a)(v_k^\top u_b) > 0$:

- 678 • if $\mu_k > 0$ (a locally convex direction), the per-direction contribution is positive, so this
 679 aligned perspective increases the second-order correction and is harmful under the loss-
 680 minimization convention;
- 681 • if $\mu_k < 0$ (a locally concave direction), the per-direction contribution is negative, so this
 682 aligned perspective decreases the second-order correction and is beneficial under the same
 683 convention.

684 The same logic applies with reversed signs when the two shifts are anti-aligned along v_k . In that case,
 685 anti-alignment along a positive-curvature direction contributes negatively, while anti-alignment along
 686 a negative-curvature direction contributes positively. The interaction term thus evaluates pairwise
 687 behavior in a curvature-aware manner. The alignment between two examples can either drive f
 688 toward worse values or pull it toward better ones, depending on the sign of the curvature along that
 689 direction.

690 **C.5 Proof of Proposition 1**

691 We prove the closed-form factorization of $\kappa(a, b)$ stated in Proposition 1. We work in the binary
 692 logistic regression setting with ℓ_2 regularization of strength $\beta > 0$, and let $\hat{\theta}$ denote the regularized
 693 empirical risk minimizer. The proof proceeds in three steps: deriving the closed-form per-example
 694 gradient and the corresponding Hessian of the regularized objective, expressing the per-example
 695 parameter shift in factored form, and substituting into the definition of $\kappa(a, b)$.

696 **Step 1: per-example gradient and training-loss Hessian.** For binary logistic regression, the
 697 per-example data loss takes the form

$$\ell(z_i, \theta) = -y_i \log \sigma(\theta^\top x_i) - (1 - y_i) \log(1 - \sigma(\theta^\top x_i)), \quad (33)$$

698 where $y_i \in \{0, 1\}$ is the binary label. Using the standard identity $\sigma'(t) = \sigma(t)(1 - \sigma(t))$ and writing
 699 $\sigma_i := \sigma(\hat{\theta}^\top x_i)$, the per-example gradient evaluated at $\hat{\theta}$ is

$$\nabla_{\theta} \ell(z_i, \hat{\theta}) = (\sigma_i - y_i) x_i. \quad (34)$$

700 The training objective is the regularized empirical risk

$$\mathcal{R}(\theta) = \frac{1}{N} \sum_i \ell(z_i, \theta) + \frac{\beta}{2} \|\theta\|^2. \quad (35)$$

701 Differentiating once gives $\nabla_{\theta} \mathcal{R}(\theta) = \frac{1}{N} \sum_i \nabla_{\theta} \ell(z_i, \theta) + \beta \theta$. Differentiating once more yields the
 702 training-loss Hessian

$$H := \nabla_{\theta}^2 \mathcal{R}(\hat{\theta}) = \frac{1}{N} \sum_i \nabla_{\theta}^2 \ell(z_i, \hat{\theta}) + \beta I = \frac{1}{N} \sum_i \sigma_i (1 - \sigma_i) x_i x_i^\top + \beta I, \quad (36)$$

703 where the last equality uses the closed-form per-example Hessian of the binary logistic loss. The data-
 704 dependent part $\frac{1}{N} \sum_i \sigma_i (1 - \sigma_i) x_i x_i^\top$ is positive semidefinite as a sum of rank-one outer products
 705 with non-negative weights, and the regularizer contributes βI with $\beta > 0$. Their sum H is therefore
 706 positive definite and hence invertible. We use Equation (34) as the per-example gradient driving the
 707 influence approximation throughout the remainder of the proof.

708 **Step 2: per-example parameter shift.** Substituting the closed-form gradient into the definition
 709 $u_i := H^{-1} \nabla_{\theta} \ell(z_i, \hat{\theta})$ gives

$$u_i = H^{-1}[(\sigma_i - y_i) x_i] = (\sigma_i - y_i) H^{-1} x_i, \quad (37)$$

710 where the second equality uses the linearity of matrix-vector multiplication and the fact that $(\sigma_i - y_i)$ is
 711 a scalar. Each per-example parameter shift therefore factors into a scalar prediction residual $(\sigma_i - y_i)$
 712 and a vector $H^{-1} x_i$ that depends only on the input feature x_i through the inverse training-loss
 713 Hessian.

714 **Step 3: factorization of $\kappa(a, b)$.** Substituting Equation (37) into the definition $\kappa(a, b) := u_a^\top H_f u_b$
 715 from Equation (9) gives

$$\begin{aligned} \kappa(a, b) &= [(\sigma_a - y_a) H^{-1} x_a]^\top H_f [(\sigma_b - y_b) H^{-1} x_b] \\ &= (\sigma_a - y_a)(\sigma_b - y_b) x_a^\top H^{-1} H_f H^{-1} x_b \\ &= (\sigma_a - y_a)(\sigma_b - y_b) \langle x_a, x_b \rangle_M, \end{aligned} \quad (38)$$

716 where the second equality factors the two scalar residuals out of the bilinear form and uses the
 717 symmetry of H^{-1} , which holds because H is symmetric and positive definite, and the third equality
 718 applies the definition $M := H^{-1} H_f H^{-1}$ together with $\langle u, v \rangle_M := u^\top M v$. This is exactly the
 719 factorization claimed in Proposition 1, completing the proof.

720 **C.6 Extension of Proposition 1 to deep classifiers**

721 Proposition 1 establishes a closed-form factorization of $\kappa(a, b)$ for binary logistic regression. We
 722 now extend this result to deep classifiers, treating binary and multi-class classification in turn. The
 723 argument is in both cases a direct consequence of chain rule applied to the per-example gradient: the
 724 same scalar-residual-input-feature factorization that drives the proof of Proposition 1 survives when
 725 the input feature x_i is replaced by the logit Jacobian

$$J_i := \nabla_{\theta} f_{\hat{\theta}}(x_i), \quad (39)$$

726 where $f_{\theta} : \mathcal{X} \rightarrow \mathbb{R}^C$ is the network logit output, with $C = 1$ for binary classification and $C \geq 2$
 727 for multi-class classification. The bilinear-form matrix $M := H^{-1} H_f H^{-1}$ from Proposition 1 is
 728 unchanged. We treat the binary case first as it offers the cleanest correspondence with Proposition 1,
 729 then state and prove the multi-class generalization.

730 **Binary case.** For binary classification ($C = 1$), $J_i \in \mathbb{R}^p$ is a vector. Writing $\sigma_i := \sigma(f_{\hat{\theta}}(x_i))$, the
 731 per-example loss $\ell(z_i, \theta) = -y_i \log \sigma(f_{\theta}(x_i)) - (1 - y_i) \log(1 - \sigma(f_{\theta}(x_i)))$ has gradient at $\hat{\theta}$ given
 732 by chain rule as

$$\nabla_{\theta} \ell(z_i, \hat{\theta}) = (\sigma_i - y_i) J_i, \quad (40)$$

733 which differs from Equation (34) only in that the input feature x_i is replaced by the logit Jacobian J_i .
 734 Steps 2 and 3 of the proof of Proposition 1 therefore carry over verbatim with $x_i \mapsto J_i$, yielding

$$\kappa(a, b) = (\sigma_a - y_a)(\sigma_b - y_b) \langle J_a, J_b \rangle_M. \quad (41)$$

735 The class-agreement scalar is unchanged because it depends only on the network’s predictions and the
 736 labels. The input-similarity term $\langle x_a, x_b \rangle_M$ generalizes to a Jacobian bilinear form $\langle J_a, J_b \rangle_M$, which
 737 is an M -weighted empirical neural tangent kernel evaluated at $\hat{\theta}$ when H_f is positive semidefinite.

738 **Decomposing the Jacobian bilinear form via the last layer.** Writing $\theta = (\psi, w)$ for the feature-
 739 extractor parameters and the last-layer weight, and decomposing $f_{\theta}(x) = w^{\top} \phi_{\psi}(x)$ where $\phi_{\psi} : \mathcal{X} \rightarrow$
 740 $\mathbb{R}^{d_{\phi}}$ denotes the learned feature representation, the logit Jacobian J_i admits the block decomposition

$$J_i = \begin{bmatrix} J_{\phi}(x_i)^{\top} w \\ \phi(x_i) \end{bmatrix}, \quad J_{\phi}(x_i) := \frac{\partial \phi_{\psi}(x_i)}{\partial \psi}, \quad (42)$$

741 corresponding to gradients with respect to ψ (top block) and w (bottom block). Substituting into
 742 the Jacobian bilinear form exposes both the learned feature representation and the feature Jacobian
 743 directly:

$$\langle J_a, J_b \rangle_M = \begin{bmatrix} J_{\phi}(x_a)^{\top} w \\ \phi(x_a) \end{bmatrix}^{\top} M \begin{bmatrix} J_{\phi}(x_b)^{\top} w \\ \phi(x_b) \end{bmatrix}. \quad (43)$$

744 In Equation (43), the bottom-block component $\phi(x_i)$ takes the place of the raw input feature x_i in
 745 Proposition 1. Since the network’s representation is well known to map visually similar examples
 746 to similar $\phi(x_i)$, such examples contribute to $\kappa(a, b)$ through this block in the same way that input
 747 similarity contributes in the LR setting. The top-block component $J_{\phi}(x_i)^{\top} w$ captures how the feature
 748 representation responds to feature-extractor parameter perturbations along the readout direction w ,
 749 an additional structure unique to deep classifiers that has no counterpart in the LR setting.

750 **Multi-class case.** For multi-class classification with $C \geq 2$, $J_i \in \mathbb{R}^{p \times C}$ is a matrix whose c -
 751 th column is the gradient of the c -th logit with respect to all parameters. Write $p_i \in \Delta^{C-1}$ for
 752 the softmax of $f_{\hat{\theta}}(x_i)$ and $y_i \in \{0, 1\}^C$ for the one-hot label. The cross-entropy loss $\ell(z_i, \theta) =$
 753 $-\sum_{c=1}^C y_{i,c} \log p_{i,c}(\theta)$ has gradient at $\hat{\theta}$ given by chain rule as

$$\nabla_{\theta} \ell(z_i, \hat{\theta}) = J_i r_i, \quad r_i := p_i - y_i \in \mathbb{R}^C. \quad (44)$$

754 Equation (44) generalizes Equation (34) in two parallel ways: the scalar residual $(\sigma_i - y_i)$ becomes
 755 the vector residual r_i , and the input feature x_i becomes the matrix Jacobian J_i .

756 Substituting Equation (44) into the definition $u_i := H^{-1} \nabla_{\theta} \ell(z_i, \hat{\theta})$ gives $u_i = H^{-1} J_i r_i$. Plugging
 757 this into the pairwise interaction $\kappa(a, b) := u_a^{\top} H_f u_b$ from Equation (9) and factoring yields

$$\begin{aligned} \kappa(a, b) &= (H^{-1} J_a r_a)^{\top} H_f (H^{-1} J_b r_b) \\ &= r_a^{\top} J_a^{\top} H^{-1} H_f H^{-1} J_b r_b \\ &= r_a^{\top} K_M(x_a, x_b) r_b, \end{aligned} \quad (45)$$

758 where the second equality uses the symmetry of H^{-1} , and the third equality defines the matrix-valued
 759 bilinear form

$$K_M(x_a, x_b) := J_a^\top M J_b \in \mathbb{R}^{C \times C}, \quad (46)$$

760 which is a matrix-valued empirical neural tangent kernel when H_f is positive semidefinite. The
 761 factorization in Equation (45) is the multi-class generalization of Equation (38): the scalar residual
 762 product $(\sigma_a - y_a)(\sigma_b - y_b)$ is replaced by the bilinear form $r_a^\top K_M r_b$, and the scalar feature
 763 similarity $\langle x_a, x_b \rangle_M$ is replaced by the matrix-valued bilinear form $K_M(x_a, x_b)$. The binary case in
 764 Equation (41) is recovered when $C = 1$, with $r_i = \sigma_i - y_i$ and $K_M = \langle J_a, J_b \rangle_M$ a scalar.

765 **Connection to the class-agreement structure.** Proposition 1 identifies cross-class similar pairs
 766 as beneficial and same-class similar pairs as redundant, with the sign of $\kappa(a, b)$ controlled by the
 767 residual product $(\sigma_a - y_a)(\sigma_b - y_b)$. The multi-class factorization in Equation (45) preserves the
 768 structural form *residual structure* \times *Jacobian similarity*, with the residual product replaced by the
 769 bilinear form $r_a^\top K_M r_b$. A direct calculation yields the residual-alignment identity

$$r_a^\top r_b = \langle p_a, p_b \rangle - p_{a, y_b} - p_{b, y_a} + \mathbb{1}[y_a = y_b], \quad (47)$$

770 in which the indicator term contributes +1 for same-class pairs and 0 for cross-class pairs. For
 771 predictions biased toward their true classes, this indicator drives the sign of the residual alignment
 772 positive for same-class pairs and negative for cross-class pairs; for example, with $p_{i,c} = (1 - \delta) \mathbb{1}[c =$
 773 $y_i] + \delta/(C - 1)$ on the remaining classes, a short calculation gives $r_a^\top r_b = \delta^2 C/(C - 1) > 0$ for
 774 same-class pairs and $r_a^\top r_b = -\delta^2 C/(C - 1)^2 < 0$ for cross-class pairs. When two examples are
 775 similar so that $J_a \approx J_b$ and M is positive semidefinite, $K_M(x_a, x_b) \approx J_a^\top M J_a$, which is symmetric
 776 PSD, and $\kappa(a, b)$ inherits the sign of the residual alignment. This recovers the cross-class-beneficial
 777 and same-class-redundant structure of Proposition 1 and is consistent with the Figure 2 observation
 778 that this class-pair structure carries over from logistic regression to MLPs and ResNet-9.

779 C.7 Derivation of the marginal score

780 We derive the marginal score in Equation (12) from the definition $m(z_j | S) := \hat{\mathcal{I}}^+(S \cup \{z_j\}) - \hat{\mathcal{I}}^+(S)$
 781 in Section 3.2. The derivation amounts to substituting the estimator $\hat{\mathcal{I}}^+(\cdot)$ into the definition and
 782 expanding the resulting quadratic form using the linearity of the aggregate parameter shift.

783 Recall from Equation (8) that the addition-setting estimator is

$$\hat{\mathcal{I}}^+(S) := -\frac{1}{N} \nabla_{\theta} f(\hat{\theta})^\top u_S + \frac{1}{2N^2} u_S^\top H_f u_S, \quad (48)$$

784 and that the aggregate approximated shift over a subset T is $u_T := \sum_{z_i \in T} u_i$. For $T = S \cup \{z_j\}$
 785 with $z_j \notin S$, the linearity of the sum gives

$$u_{S \cup \{z_j\}} = u_S + u_j. \quad (49)$$

786 Substituting Equation (49) into the linear term of Equation (48) evaluated at $S \cup \{z_j\}$ gives

$$-\frac{1}{N} \nabla_{\theta} f(\hat{\theta})^\top u_{S \cup \{z_j\}} = -\frac{1}{N} \nabla_{\theta} f(\hat{\theta})^\top u_S - \frac{1}{N} \nabla_{\theta} f(\hat{\theta})^\top u_j. \quad (50)$$

787 The first term on the right cancels with the corresponding linear term in $\hat{\mathcal{I}}^+(S)$, leaving only the
 788 contribution of z_j .

789 Substituting Equation (49) into the quadratic term and expanding gives

$$\begin{aligned} \frac{1}{2N^2} u_{S \cup \{z_j\}}^\top H_f u_{S \cup \{z_j\}} &= \frac{1}{2N^2} (u_S + u_j)^\top H_f (u_S + u_j) \\ &= \frac{1}{2N^2} u_S^\top H_f u_S + \frac{1}{N^2} u_S^\top H_f u_j + \frac{1}{2N^2} u_j^\top H_f u_j, \end{aligned} \quad (51)$$

790 where the cross terms $u_S^\top H_f u_j$ and $u_j^\top H_f u_S$ combine into a single term with coefficient $1/N^2$
 791 because H_f is symmetric, which gives $u_j^\top H_f u_S = u_S^\top H_f u_j$. The first term on the right of
 792 Equation (51) is the quadratic term in $\hat{\mathcal{I}}^+(S)$ itself and cancels in the difference $\hat{\mathcal{I}}^+(S \cup \{z_j\}) -$
 793 $\hat{\mathcal{I}}^+(S)$.

794 Combining the surviving contributions from Equations (50) and (51) yields

$$m(z_j | S) \approx -\frac{1}{N} \nabla_{\theta} f(\hat{\theta})^{\top} u_j + \frac{1}{N^2} u_S^{\top} H_f u_j + \frac{1}{2N^2} u_j^{\top} H_f u_j, \quad (52)$$

795 which is exactly Equation (12) in the main text. The three surviving terms admit the interpretation
796 given in Section 3.2: the first term is the standard first-order influence of the candidate z_j , the second
797 term measures the curvature-driven interaction between z_j and the already-selected subset S , and the
798 third term is the candidate’s self-contribution to the curvature correction.

799 This decomposition aligns with the self versus cross distinction developed in Appendix C.3. The
800 cross term $\frac{1}{N^2} u_S^{\top} H_f u_j$ in Equation (52) is precisely the contribution of z_j to the cross part of the
801 pairwise decomposition, while the self term $\frac{1}{2N^2} u_j^{\top} H_f u_j$ is its contribution to the diagonal part.
802 The greedy procedure in Algorithm 1 therefore exploits the same self versus cross structure to update
803 each candidate’s marginal utility based on its interaction with previously selected examples, without
804 recomputing the full quadratic form at each iteration.

805 D Scalable Approximations: EK-FAC and Low-Rank Gradient Projection

806 The estimators $\hat{\mathcal{L}}^-(S)$ and $\hat{\mathcal{L}}^+(S)$ defined in Equation (7) and Equation (8) require applying H^{-1}
807 to per-example gradients. Since H has p^2 entries in p parameters and is moreover not guaranteed
808 to be positive definite away from a minimum, this is intractable for the models considered in our
809 experiments. We therefore approximate H by the damped Gauss–Newton matrix $G + \lambda I$, where G is
810 the Gauss–Newton Hessian of the training loss at $\hat{\theta}$ and $\lambda > 0$ is a damping constant. This appendix
811 collects the three approximations that make our estimator tractable at scale: the Gauss–Newton
812 approximation itself, its block-diagonal EK-FAC factorization, and a low-rank gradient projection for
813 large language models. The damping constant λ is treated as a hyperparameter and its selection is
814 described in Appendix E.

815 D.1 Damped Gauss–Newton approximation

816 Following standard practice in influence-function analysis [4, 22], we approximate the training
817 Hessian by the Gauss–Newton Hessian G of the training loss at $\hat{\theta}$, damped by λI to ensure positive
818 definiteness. The Gauss–Newton Hessian is positive semidefinite by construction and discards the
819 indefinite residual term that arises in the exact Hessian, making it a natural surrogate when the goal is
820 to invert a curvature matrix.

821 The losses considered in this paper are all negative log-likelihoods of exponential-family distributions:
822 cross-entropy with softmax outputs for classification, squared error with a Gaussian likelihood
823 for regression, and token-level negative log-likelihood for instruction tuning. For such losses, the
824 Gauss–Newton Hessian coincides with the Fisher information matrix taken under the model’s output
825 distribution [37]. This equivalence is what licenses the use of curvature approximations developed
826 for the Fisher information matrix, and it is the same identification adopted by Grosse et al. [22] for
827 influence-function analysis of transformer language models. We accordingly state our approximations
828 in terms of G throughout, with the understanding that G may equivalently be read as the model-
829 distribution Fisher information.

830 D.2 Block-diagonal structure and EK-FAC

831 Even with the Gauss–Newton substitution, G remains a $p \times p$ matrix and cannot be inverted directly
832 at the parameter counts considered in our experiments. We therefore further approximate G as
833 block-diagonal across layers, treating cross-layer curvature as zero. This reduces the inversion to one
834 block per layer.

835 Within each block, we apply the Eigenvalue-Corrected Kronecker-Factored Approximate Curvature
836 (EK-FAC) approximation [19, 22]. EK-FAC builds on K-FAC [38], which factorizes the per-layer
837 Gauss–Newton block as a Kronecker product of two small matrices: the second-moment matrix of
838 the layer’s input activations and the second-moment matrix of the gradient with respect to the layer’s
839 pre-activations. EK-FAC retains the eigenbasis defined by these Kronecker factors but replaces the
840 implied Kronecker-structured eigenvalues with the exact diagonal of G in that eigenbasis, estimated

841 from training data. The resulting approximation is provably at least as accurate as K-FAC under
 842 the Frobenius norm [19] while preserving the same favorable scaling: each per-layer block can be
 843 inverted, and inverse-vector products evaluated, at a cost dominated by the layer’s activation and
 844 pre-activation dimensions rather than its full parameter count. In practice, the per-example shifts
 845 u_i required by Equation (7) and Equation (8) are computed as $(G + \lambda I)^{-1} \nabla_{\theta} \ell(z_i, \hat{\theta})$ by applying
 846 this block-wise EK-FAC inverse to per-example gradients. We refer the reader to George et al. [19]
 847 for the original derivation of EK-FAC and to Grosse et al. [22] for a detailed treatment of its use in
 848 influence-function analysis.

849 D.3 Low-rank gradient projection for large language models

850 At the scale of Llama-3.1-8B, the per-layer EK-FAC blocks remain large enough that materializing
 851 and storing one projected gradient u_i per training example, as required by the greedy procedure in
 852 Algorithm 1, is itself the dominant cost. To address this, we adopt the low-rank gradient projection of
 853 LoGra [10], which introduces a shared low-dimensional subspace and projects all per-example gradi-
 854 ents into that subspace before any subsequent inner product or curvature operation is performed. The
 855 projection is applied consistently to gradients and to the target-side quantities entering Algorithm 1,
 856 so that all inner products in the greedy loop are evaluated in the projected dimension d rather than in
 857 the full parameter dimension p .

858 This projection is what reduces the per-iteration cost of Algorithm 1 to $O(P)$ inner products in
 859 dimension d , yielding the overall greedy-loop complexity of $O(KPd)$ stated in Section 3.2. We
 860 use the projection dimension and hyperparameters reported in Appendix E, and refer the reader to
 861 Choe et al. [10] for the construction of the projection operator and a full treatment of the resulting
 862 estimator.

863 E Experimental Details

864 E.1 Small-scale Attribution Experiments

865 The attribution experiments in Section 4 consider six dataset–model pairs: logistic regression and
 866 MLP on MNIST, MLP on FashionMNIST, MLP on Concrete and Parkinsons, and ResNet-9 on
 867 CIFAR-10. All datasets use their standard train/test splits.

868 **Models.** LR uses ℓ_2 regularization with strength 0.01. The MLP has two hidden layers of width
 869 128 and 64 with ReLU activations. ResNet-9 follows the standard DAWNBench architecture [14],
 870 trained from scratch.

871 **Training.** Training configurations are summarized in Table 3 and follow standard settings for each
 model.

Table 3: Training configurations for small-scale models. The cosine schedule decays the learning rate to zero by the end of training.

Model	Optimizer	Momentum	LR	LR Schedule	Weight Decay	Batch Size	Epochs
LR	SGD	0.0	0.01	constant	0.01	64	200
MLP	SGD	0.0	0.01	constant	0.01	64	200
ResNet-9	SGD	0.9	0.01	cosine	0.01	64	50

872

873 **Group construction.** Each group consists of an anchor sampled uniformly at random from the
 874 training set, together with its $|S| - 1$ nearest neighbors in softmax output space (under L^2 distance).
 875 We set $|S| = 400$, yielding 50 subsets per dataset–model pair, except for Concrete where we set
 876 $|S| = 100$ due to its small training set. For regression tasks, the prediction itself replaces the softmax
 877 output.

878 **Ground truth.** For each subset S , we retrain the model from scratch on $\mathcal{D} \setminus S$ and report the mean
 879 change in held-out test loss.

Table 4: Compute cost of selection methods on the LESS pool with Llama-3.1-8B, in GPU-hours.

Method	Warmup	Selection	Total
Additive IF	9.2	22	31
LESS [59]	9.2	306	315
RDS+ [27]	–	23	23
NV-Embed [35]	–	109	109
Ours	9.2	22	31

Table 5: Per-component GPU-hours of our selection pipeline on the LESS pool with Llama-3.1-8B.

Component	GPU-hours	Scope
Warmup fine-tuning	9.2	shared
Pool gradient extraction (LoGra)	22	shared
Selection step	0.07	per target

880 **Influence computation.** For LR, the Hessian is computed and inverted exactly. For MLPs and
 881 ResNet-9, we use the damped Gauss–Newton approximation of the training Hessian, $G + \lambda I$, and
 882 approximate its inverse via EK-FAC. The damping λ is set to 10^{-2} for all influence-based methods.
 883 The target Hessian H_f is approximated as block-diagonal, with H_f -vector products computed directly
 884 without inversion.

885 **Baselines.** F is the standard first-order influence function [29]. $F+B$ applies the second-order
 886 response-function correction of Basu et al. [5], computed under the same curvature approximation
 887 as our method. $F+I+B$ combines our interaction term with the Basu correction under identical
 888 approximations. For TRAK [43], we use projection dimension 2048 and ridge 0.01. For TracIn [44],
 889 we use 10 evenly spaced checkpoints. For both TRAK and TracIn, group scores are obtained by
 890 summing individual attributions.

891 E.2 Small-Scale Selection Experiments

892 The selection experiments in Section 5.1 use a two-hidden-layer MLP ($784 \rightarrow 128 \rightarrow 64 \rightarrow 10$) with
 893 ReLU activations and dropout rate 0.1. For each dataset, a candidate pool $\mathcal{D}_{\text{pool}}$ of 5,000 training
 894 examples is drawn by a fixed permutation of the original training set, and the reference parameter $\hat{\theta}$
 895 is obtained by training the MLP on this pool for 200 epochs with vanilla SGD (learning rate 10^{-2} ,
 896 batch size 64, weight decay 3×10^{-2}). The Hessian-inverse H^{-1} is approximated via EK-FAC with
 897 the GGN Fisher and damping $\lambda = 5 \times 10^{-1}$. Each selection picks $K \in \{500, 1,000, \dots, 5,000\}$
 898 examples in increments of 500 from the same pool $\mathcal{D}_{\text{pool}}$, and each selected subset is retrained from
 899 scratch with the same training configuration. Held-out test loss is averaged over five random seeds
 900 and reported; the class-entropy panel reports the Shannon entropy of the empirical class distribution
 901 among the K selected examples.

902 E.3 LLM Data Selection Experiments

903 We fine-tune Llama-3.1-8B on subsets selected from the LESS pool of approximately 270K
 904 instruction-tuning examples.

905 **Warmup and selection.** Following Xia et al. [59], a warmup model $\hat{\theta}$ is obtained by fine-tuning the
 906 base model on a random subset of size 13,534 using LoRA (rank 8, scaling 16, applied to q_proj,
 907 k_proj, v_proj, o_proj, gate_proj, up_proj, and down_proj), AdamW with peak learning
 908 rate 10^{-4} , cosine schedule with warmup ratio 0.03, effective batch size 128, sequence length 2048,
 909 in bf16, for four epochs. At $\hat{\theta}$, we compute per-example projected gradients u_j and curvature terms
 910 $H_f u_j$ using LoGra [10] with projection dimension 16, and apply Algorithm 1 to select $K = 13,534$
 911 examples (5% of the pool) independently for each target task. The target function f is the instruction-
 912 masked likelihood, evaluated on the official task validation set when available and on a held-out
 913 training split otherwise. The selected indices for all seven target tasks are released alongside our code,
 914 enabling reproduction of the fine-tuning results without rerunning the selection pipeline.

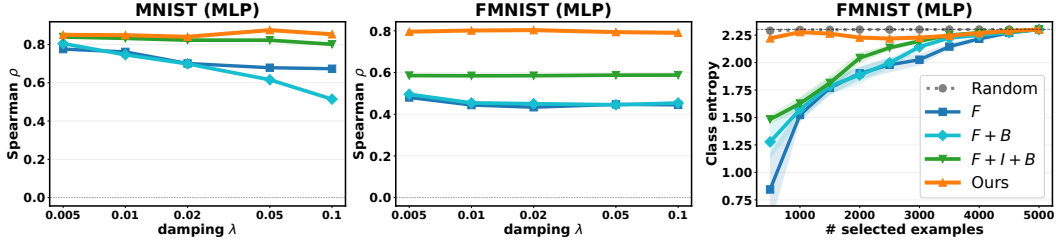


Figure 4: Sensitivity analysis of the two-layer MLP over the damping hyperparameter on MNIST (left) and FashionMNIST (middle), and class entropy on FashionMNIST (right).

915 **Final fine-tuning.** Selected subsets are used to fine-tune the original pretrained Llama-3.1-8B
 916 checkpoint (rather than the warmup model) under the same LoRA configuration as warmup: AdamW
 917 with peak learning rate 10^{-4} , cosine schedule with warmup ratio 0.03, effective batch size 128,
 918 sequence length 2048, and bf16 precision. We train for four epochs on GSM8K, whose loss does not
 919 converge within a single epoch, and one epoch on the remaining tasks.

920 **Evaluation.** We evaluate on GSM8K, AQuA, ARC-Easy, HellaSwag, PIQA, ECQA, and SQuAD
 921 using the Tulu chat template [57]. All tasks use 0-shot prompting. We report exact match for GSM8K,
 922 F1 for SQuAD, and accuracy for the multiple-choice tasks, with the latter scored by length-normalized
 923 log-likelihood. Results are averaged over five random seeds.

924 **Baselines.** All methods share the same candidate pool, selection budget, and final fine-tuning
 925 protocol; influence-based methods additionally share the same warmup checkpoint. Additive IF
 926 uses our warmup checkpoint and projected gradients but selects the top- K examples by first-order
 927 influence, omitting the interaction term. LESS [59] uses Adam-adjusted gradient cosine similarity
 928 with projection dimension 8192 on the same warmup checkpoint. RDS+ [27] uses position-weighted
 929 mean-pooled hidden states from the warmup checkpoint with cosine similarity. NV-Embed [35] uses
 930 embeddings from `nvidia/NV-Embed-v2`. Random uniformly subsamples from the pool.

931 E.4 Compute Resources

932 All small-scale experiments were conducted on a single NVIDIA RTX A6000 (48 GB), and all LLM
 933 experiments on $8 \times$ NVIDIA RTX A6000 (48 GB).

934 Table 4 reports the GPU-hours required by each selection method on the LESS pool of approximately
 935 270K instruction-tuning examples, summed over all devices used. Warmup fine-tunes a LoRA adapter
 936 on a random subset of size 13,534 for four epochs. Additive IF and our method share the same
 937 warmup checkpoint and LoGra-projected gradients, so their selection costs differ only in the greedy
 938 step, which contributes approximately 50 seconds per target.

939 Table 5 decomposes the selection cost of our pipeline into its individual components. Warmup
 940 fine-tuning and pool gradient extraction are shared across all seven targets, while the selection step
 941 is paid once per target. The selection step itself further decomposes into four sub-steps: a one-time
 942 covariance-state merge that is shared across targets, validation-set gradient extraction that scales per
 943 target, first-order influence and target-curvature inner products, and the greedy selection loop. These
 944 four sub-steps together account for the 0.07 GPU-hours per target reported in Table 5.

945 F Additional Experiments

946 In this section, we present two additional experiments on the two-layer MLP that complement the
 947 small-scale attribution accuracy results of Section 4 and the small-scale selection quality results of
 948 Section 5.1. The first examines the sensitivity of each method to the damping hyperparameter λ used
 949 in the EK-FAC, and the second extends the class entropy analysis from MNIST to FashionMNIST.
 950 The results are summarized in Figure 4.

951 **Setup.** Both experiments inherit the protocols of their counterparts in the main text, modified as
952 follows. For the damping sensitivity analysis, we follow the attribution setup of Section 4 but vary
953 $\lambda \in \{0.005, 0.01, 0.02, 0.05, 0.1\}$ instead of fixing $\lambda = 10^{-2}$. All other elements of the protocol,
954 including group construction, ground-truth retraining, and the influence-function methods (F , $F+B$,
955 $F+I+B$, Ours), are unchanged. For the class entropy analysis, we run the selection procedure of
956 Section 5.1 on FashionMNIST instead of MNIST, keeping all other choices unchanged.

957 **Damping sensitivity.** The left and middle panels of Figure 4 report the Spearman correlation as a
958 function of the damping λ on MNIST and FashionMNIST. Ours and F+I+B remain nearly constant
959 across the tested range on both datasets, while F and F+B are somewhat less stable, with F+B showing
960 the most pronounced decline at larger λ on MNIST.

961 **Class entropy on FashionMNIST.** The right panel of Figure 4 reports the class entropy of the
962 subsets selected by each method on FashionMNIST. The pattern matches the MNIST result reported
963 in Figure 3: at small budgets, F and F+B collapse onto a few classes and F+I+B shows a milder but
964 still noticeable drop, while our method matches the entropy of random selection across all selection
965 sizes. This confirms that the diversity-promoting behavior of the interaction term is not specific to
966 MNIST and carries over to FashionMNIST.

967 **NeurIPS Paper Checklist**

968 **1. Claims**

969 Question: Do the main claims made in the abstract and introduction accurately reflect the
970 paper’s contributions and scope?

971 Answer: [Yes]

972 Justification: The abstract and introduction state three claims: (i) we propose an interaction-
973 aware influence function from a second-order target-side expansion (formalized in Section 3),
974 (ii) for logistic regression, the pairwise interaction term identifies cross-class pairs as
975 beneficial and same-class pairs as redundant (Proposition 1), and (iii) the estimator improves
976 Spearman correlation with retraining ground truth by up to 0.67 across six dataset–model
977 pairs and outperforms random selection on Llama-3.1-8B instruction tuning (Section 5).
978 Each claim is supported by the corresponding theoretical or empirical result.

979 Guidelines:

- 980 • The answer [N/A] means that the abstract and introduction do not include the claims
981 made in the paper.
- 982 • The abstract and/or introduction should clearly state the claims made, including the
983 contributions made in the paper and important assumptions and limitations. A [No] or
984 [N/A] answer to this question will not be perceived well by the reviewers.
- 985 • The claims made should match theoretical and experimental results, and reflect how
986 much the results can be expected to generalize to other settings.
- 987 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
988 are not attained by the paper.

989 **2. Limitations**

990 Question: Does the paper discuss the limitations of the work performed by the authors?

991 Answer: [Yes]

992 Justification: We discuss limitations in the conclusion, noting that our greedy selection
993 procedure lacks formal optimality guarantees and that our estimator inherits approximation
994 error from the underlying Hessian approximations.

995 Guidelines:

- 996 • The answer [N/A] means that the paper has no limitation while the answer [No] means
997 that the paper has limitations, but those are not discussed in the paper.
- 998 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 999 • The paper should point out any strong assumptions and how robust the results are to
1000 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1001 model well-specification, asymptotic approximations only holding locally). The authors
1002 should reflect on how these assumptions might be violated in practice and what the
1003 implications would be.
- 1004 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1005 only tested on a few datasets or with a few runs. In general, empirical results often
1006 depend on implicit assumptions, which should be articulated.
- 1007 • The authors should reflect on the factors that influence the performance of the approach.
1008 For example, a facial recognition algorithm may perform poorly when image resolution
1009 is low or images are taken in low lighting. Or a speech-to-text system might not be
1010 used reliably to provide closed captions for online lectures because it fails to handle
1011 technical jargon.
- 1012 • The authors should discuss the computational efficiency of the proposed algorithms
1013 and how they scale with dataset size.
- 1014 • If applicable, the authors should discuss possible limitations of their approach to
1015 address problems of privacy and fairness.
- 1016 • While the authors might fear that complete honesty about limitations might be used by
1017 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1018 limitations that aren’t acknowledged in the paper. The authors should use their best

1019 judgment and recognize that individual actions in favor of transparency play an impor-
1020 tant role in developing norms that preserve the integrity of the community. Reviewers
1021 will be specifically instructed to not penalize honesty concerning limitations.

1022 3. Theory assumptions and proofs

1023 Question: For each theoretical result, does the paper provide the full set of assumptions and
1024 a complete (and correct) proof?

1025 Answer: [Yes]

1026 Justification: The first-order influence approximation states the standard assumptions (twice
1027 continuous differentiability of ℓ , nonsingularity of H) in Section 2, and the derivation is
1028 provided in Appendix B. Proposition 1 states its assumption (binary logistic regression with
1029 ℓ_2 regularization) and is proved in Appendix C.5. All theorems and equations are numbered
1030 and cross-referenced.

1031 Guidelines:

- 1032 • The answer [N/A] means that the paper does not include theoretical results.
- 1033 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
1034 referenced.
- 1035 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1036 • The proofs can either appear in the main paper or the supplemental material, but if
1037 they appear in the supplemental material, the authors are encouraged to provide a short
1038 proof sketch to provide intuition.
- 1039 • Inversely, any informal proof provided in the core of the paper should be complemented
1040 by formal proofs provided in appendix or supplemental material.
- 1041 • Theorems and Lemmas that the proof relies upon should be properly referenced.

1042 4. Experimental result reproducibility

1043 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
1044 perimental results of the paper to the extent that it affects the main claims and/or conclusions
1045 of the paper (regardless of whether the code and data are provided or not)?

1046 Answer: [Yes]

1047 Justification: Algorithm 1 specifies the selection procedure end-to-end. Datasets, model
1048 architectures, training configurations, group construction, ground-truth retraining protocol,
1049 EK-FAC and damping settings, baseline hyperparameters, warmup pipeline, LoRA configu-
1050 ration, evaluation protocol, and seed averaging are described in Appendix E. All datasets
1051 used are publicly available.

1052 Guidelines:

- 1053 • The answer [N/A] means that the paper does not include experiments.
- 1054 • If the paper includes experiments, a [No] answer to this question will not be perceived
1055 well by the reviewers: Making the paper reproducible is important, regardless of
1056 whether the code and data are provided or not.
- 1057 • If the contribution is a dataset and/or model, the authors should describe the steps taken
1058 to make their results reproducible or verifiable.
- 1059 • Depending on the contribution, reproducibility can be accomplished in various ways.
1060 For example, if the contribution is a novel architecture, describing the architecture fully
1061 might suffice, or if the contribution is a specific model and empirical evaluation, it may
1062 be necessary to either make it possible for others to replicate the model with the same
1063 dataset, or provide access to the model. In general, releasing code and data is often
1064 one good way to accomplish this, but reproducibility can also be provided via detailed
1065 instructions for how to replicate the results, access to a hosted model (e.g., in the case
1066 of a large language model), releasing of a model checkpoint, or other means that are
1067 appropriate to the research performed.
- 1068 • While NeurIPS does not require releasing code, the conference does require all submis-
1069 sions to provide some reasonable avenue for reproducibility, which may depend on the
1070 nature of the contribution. For example
1071 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
1072 to reproduce that algorithm.

- 1073 (b) If the contribution is primarily a new model architecture, the paper should describe
1074 the architecture clearly and fully.
- 1075 (c) If the contribution is a new model (e.g., a large language model), then there should
1076 either be a way to access this model for reproducing the results or a way to reproduce
1077 the model (e.g., with an open-source dataset or instructions for how to construct
1078 the dataset).
- 1079 (d) We recognize that reproducibility may be tricky in some cases, in which case
1080 authors are welcome to describe the particular way they provide for reproducibility.
1081 In the case of closed-source models, it may be that access to the model is limited in
1082 some way (e.g., to registered users), but it should be possible for other researchers
1083 to have some path to reproducing or verifying the results.

1084 5. Open access to data and code

1085 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1086 tions to faithfully reproduce the main experimental results, as described in supplemental
1087 material?

1088 Answer: [Yes]

1089 Justification: An anonymized code repository containing implementations of our estimator,
1090 the greedy selection procedure, the selected indices for all seven target tasks, and scripts
1091 to reproduce the small-scale and LLM experiments is provided as supplementary material
1092 at https://anonymous.4open.science/r/Interaction_IF-45D6. All datasets used
1093 (MNIST, FashionMNIST, Concrete, Parkinsons, CIFAR-10, the LESS instruction-tuning
1094 pool, and the seven evaluation tasks) are publicly available; instructions for accessing and
1095 preprocessing them are included in the repository README.

1096 Guidelines:

- 1097 • The answer [N/A] means that paper does not include experiments requiring code.
- 1098 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
1099 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1100 • While we encourage the release of code and data, we understand that this might not
1101 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
1102 including code, unless this is central to the contribution (e.g., for a new open-source
1103 benchmark).
- 1104 • The instructions should contain the exact command and environment needed to run to
1105 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1106 • The authors should provide instructions on data access and preparation, including how
1107 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1108 • The authors should provide scripts to reproduce all experimental results for the new
1109 proposed method and baselines. If only a subset of experiments are reproducible, they
1110 should state which ones are omitted from the script and why.
- 1111 • At submission time, to preserve anonymity, the authors should release anonymized
1112 versions (if applicable).
- 1113 • Providing as much information as possible in supplemental material (appended to the
1114 paper) is recommended, but including URLs to data and code is permitted.
- 1115

1116 6. Experimental setting/details

1117 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
1118 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1119 Answer: [Yes]

1120 Justification: For the small-scale attribution experiments, optimizer, momentum, learning
1121 rate and schedule, weight decay, batch size, and training budget for each model are sum-
1122 marized in Table 3; group construction, damping, and baseline-specific hyperparameters
1123 are described in Appendix E. For the LLM data selection experiments, warmup and final
1124 fine-tuning configurations (LoRA setup, AdamW, scheduler, batch size, sequence length,
1125 precision, epochs) and selection-pipeline hyperparameters (LoGra projection dimension,
1126 selection budget) are also provided in Appendix E.

1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Table 1 reports mean \pm standard deviation over 5 random seeds for the LLM experiments, where the seed varies fine-tuning initialization and data shuffling.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Compute resources, including GPU type, number of devices, and approximate GPU-hours per fine-tuning run and total, are reported in the Compute Resources subsection of Appendix E.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

1178 Question: Does the research conducted in the paper conform, in every respect, with the
1179 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1180 Answer: [Yes]

1181 Justification: The research uses publicly available datasets and pretrained models, does not
1182 involve human subjects, and conforms with the NeurIPS Code of Ethics.

1183 Guidelines:

- 1184 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
1185 Ethics.
- 1186 • If the authors answer [No], they should explain the special circumstances that require a
1187 deviation from the Code of Ethics.
- 1188 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1189 eration due to laws or regulations in their jurisdiction).

1190 10. Broader impacts

1191 Question: Does the paper discuss both potential positive societal impacts and negative
1192 societal impacts of the work performed?

1193 Answer: [N/A]

1194 Justification: This work develops a foundational data attribution and selection method. It
1195 is not tied to a specific deployment or application, and we do not foresee a direct path to
1196 negative societal impact beyond the standard considerations that apply to any data-attribution
1197 or instruction-tuning methodology.

1198 Guidelines:

- 1199 • The answer [N/A] means that there is no societal impact of the work performed.
- 1200 • If the authors answer [N/A] or [No], they should explain why their work has no societal
1201 impact or why the paper does not address societal impact.
- 1202 • Examples of negative societal impacts include potential malicious or unintended uses
1203 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1204 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1205 groups), privacy considerations, and security considerations.
- 1206 • The conference expects that many papers will be foundational research and not tied
1207 to particular applications, let alone deployments. However, if there is a direct path to
1208 any negative applications, the authors should point it out. For example, it is legitimate
1209 to point out that an improvement in the quality of generative models could be used to
1210 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
1211 that a generic algorithm for optimizing neural networks could enable people to train
1212 models that generate Deepfakes faster.
- 1213 • The authors should consider possible harms that could arise when the technology is
1214 being used as intended and functioning correctly, harms that could arise when the
1215 technology is being used as intended but gives incorrect results, and harms following
1216 from (intentional or unintentional) misuse of the technology.
- 1217 • If there are negative societal impacts, the authors could also discuss possible mitigation
1218 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1219 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1220 feedback over time, improving the efficiency and accessibility of ML).

1221 11. Safeguards

1222 Question: Does the paper describe safeguards that have been put in place for responsible
1223 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1224 image generators, or scraped datasets)?

1225 Answer: [N/A]

1226 Justification: We do not release new pretrained models, generative models, or scraped
1227 datasets. Our LLM experiments fine-tune the publicly released Llama-3.1-8B model on a
1228 publicly available instruction-tuning pool.

1229 Guidelines:

- 1230 • The answer [N/A] means that the paper poses no such risks.
- 1231 • Released models that have a high risk for misuse or dual-use should be released with
- 1232 necessary safeguards to allow for controlled use of the model, for example by requiring
- 1233 that users adhere to usage guidelines or restrictions to access the model or implementing
- 1234 safety filters.
- 1235 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 1236 should describe how they avoided releasing unsafe images.
- 1237 • We recognize that providing effective safeguards is challenging, and many papers do
- 1238 not require this, but we encourage authors to take this into account and make a best
- 1239 faith effort.

1240 12. Licenses for existing assets

1241 Question: Are the creators or original owners of assets (e.g., code, data, models), used in

1242 the paper, properly credited and are the license and terms of use explicitly mentioned and

1243 properly respected?

1244 Answer: [Yes]

1245 Justification: All datasets (MNIST, FashionMNIST, Concrete, Parkinsons, CIFAR-10,

1246 GSM8K, AQuA, ARC-Easy, HellaSwag, PIQA, ECQA, SQuAD), the LESS instruction-

1247 tuning pool, the Llama-3.1-8B model (used under the Llama 3.1 Community License), and

1248 external tools (EK-FAC, LoGra, Tulu chat template) are cited at first use. We use them under

1249 their respective public licenses and terms of use.

1250 Guidelines:

- 1251 • The answer [N/A] means that the paper does not use existing assets.
- 1252 • The authors should cite the original paper that produced the code package or dataset.
- 1253 • The authors should state which version of the asset is used and, if possible, include a
- 1254 URL.
- 1255 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1256 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 1257 service of that source should be provided.
- 1258 • If assets are released, the license, copyright information, and terms of use in the
- 1259 package should be provided. For popular datasets, paperswithcode.com/datasets
- 1260 has curated licenses for some datasets. Their licensing guide can help determine the
- 1261 license of a dataset.
- 1262 • For existing datasets that are re-packaged, both the original license and the license of
- 1263 the derived asset (if it has changed) should be provided.
- 1264 • If this information is not available online, the authors are encouraged to reach out to
- 1265 the asset's creators.

1266 13. New assets

1267 Question: Are new assets introduced in the paper well documented and is the documentation

1268 provided alongside the assets?

1269 Answer: [Yes]

1270 Justification: We release an anonymized code repository at [https://anonymous.4open.](https://anonymous.4open.science/r/Interaction_IF-45D6)

1271 [science/r/Interaction_IF-45D6](https://anonymous.4open.science/r/Interaction_IF-45D6) containing implementations of our interaction-aware

1272 influence estimator, the greedy selection procedure, the selected indices for all seven target

1273 tasks, and scripts to reproduce the small-scale attribution and LLM data selection experi-

1274 ments. The repository includes a README describing dependencies, dataset access and

1275 preprocessing, and step-by-step instructions to reproduce each experiment. We do not

1276 release new datasets or pretrained models.

1277 Guidelines:

- 1278 • The answer [N/A] means that the paper does not release new assets.
- 1279 • Researchers should communicate the details of the dataset/code/model as part of their
- 1280 submissions via structured templates. This includes details about training, license,
- 1281 limitations, etc.

- 1282 • The paper should discuss whether and how consent was obtained from people whose
1283 asset is used.
- 1284 • At submission time, remember to anonymize your assets (if applicable). You can either
1285 create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

1287 Question: For crowdsourcing experiments and research with human subjects, does the paper
1288 include the full text of instructions given to participants and screenshots, if applicable, as
1289 well as details about compensation (if any)?

1290 Answer: [N/A]

1291 Justification: The paper does not involve crowdsourcing or research with human subjects.

1292 Guidelines:

- 1293 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1294 with human subjects.
- 1295 • Including this information in the supplemental material is fine, but if the main contribu-
1296 tion of the paper involves human subjects, then as much detail as possible should be
1297 included in the main paper.
- 1298 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1299 or other labor should be paid at least the minimum wage in the country of the data
1300 collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

1303 Question: Does the paper describe potential risks incurred by study participants, whether
1304 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1305 approvals (or an equivalent approval/review based on the requirements of your country or
1306 institution) were obtained?

1307 Answer: [N/A]

1308 Justification: The paper does not involve research with human subjects and therefore does
1309 not require IRB approval.

1310 Guidelines:

- 1311 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1312 with human subjects.
- 1313 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1314 may be required for any human subjects research. If you obtained IRB approval, you
1315 should clearly state this in the paper.
- 1316 • We recognize that the procedures for this may vary significantly between institutions
1317 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1318 guidelines for their institution.
- 1319 • For initial submissions, do not include any information that would break anonymity (if
1320 applicable), such as the institution conducting the review.

16. Declaration of LLM usage

1322 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1323 non-standard component of the core methods in this research? Note that if the LLM is used
1324 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
1325 scientific rigor, or originality of the research, declaration is not required.

1326 Answer: [N/A]

1327 Justification: LLMs are not used as a component of the core methodology. Llama-3.1-8B
1328 appears as the experimental subject for instruction-tuning data selection only.

1329 Guidelines:

- 1330 • The answer [N/A] means that the core method development in this research does not
1331 involve LLMs as any important, original, or non-standard components.
- 1332 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
1333 be described.