

Consolidation or Adaptation? PRISM: Disentangling SFT and RL Data via Gradient Concentration

Anonymous ACL submission

Abstract

While Hybrid Supervised Fine-Tuning (SFT) followed by Reinforcement Learning (RL) has become the standard paradigm for training LLM agents, effective mechanisms for data allocation between these stages remain largely underexplored. Current data arbitration strategies often rely on surface-level heuristics that fail to diagnose intrinsic learning needs. Since SFT targets pattern consolidation through imitation while RL drives structural adaptation via exploration, misaligning data with these functional roles causes severe optimization interference. We propose PRISM, a dynamics-aware framework grounded in Schema Theory that arbitrates data based on its degree of cognitive conflict with the model’s existing knowledge. By analyzing the spatial geometric structure of gradients, PRISM identifies data triggering high spatial concentration as high-conflict signals that require RL for structural restructuring. In contrast, data yielding diffuse updates is routed to SFT for efficient consolidation. Extensive experiments on WebShop and ALFWorld demonstrate that PRISM achieves a Pareto improvement, outperforming state-of-the-art hybrid methods while reducing computational costs by up to 3.22 \times . Our findings suggest that disentangling data based on internal optimization regimes is crucial for scalable and robust agent alignment.

1 Introduction

Large Language Model (LLM) agents have demonstrated remarkable capabilities in complex decision-making tasks (Qian et al., 2025; Qin et al., 2024). To unlock these potentials, the prevailing paradigm has adopted a standard two-stage training pipeline: SFT to establish behavioral norms, followed by RL to optimize long-horizon exploration (Zhang et al., 2025). This pipeline relies on a functional synergy: SFT facilitates pattern consolidation by internalizing behavioral norms and task-specific

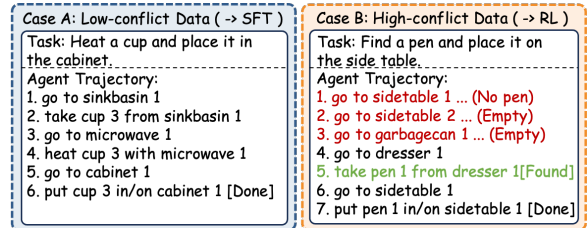


Figure 1: Case Study on ALFWorld. PRISM performs data arbitration by diagnosing cognitive conflict between task trajectories and the model’s internal state. (Left) Case A: A routine task follows a linear execution sequence, characterized by diffuse gradient updates (low concentration). Such samples are routed to SFT for behavioral consolidation. (Right) Case B: A high-conflict task involving extensive trial-and-error (e.g., searching multiple locations) triggers concentrated gradient updates (high concentration). These signals indicate a failure in the model’s current logic, necessitating RL for structural adaptation and reasoning refinement.

knowledge, while RL drives structural adaptation via trial-and-error to refine logic and enhance generalization (Chu et al., 2025). Given the divergent optimization mechanisms of SFT (imitation) and RL (exploration), the alignment efficiency and effectiveness hinges on precisely partitioning data between these regimes according to their intrinsic cognitive demands (Lv et al., 2025).

However, efficiently partitioning data between SFT and RL remains a non-trivial challenge, as current paradigms are often constrained by three primary limitations: (1) Monolithic Sequencing applies a fixed SFT-then-RL schedule to large instruction corpora (Ouyang et al., 2022). This uniform approach ignores data heterogeneity, leading to computational inefficiencies by failing to distinguish between data requiring behavioral consolidation and that necessitating exploratory reasoning (Zhou et al., 2023). (2) Universal Exploration (Shao et al., 2024; Feng et al., 2025) subjects broad trajectories to RL indiscriminately. Yet, ap-

plying trial-and-error to high-conflict data without SFT-consolidated behavioral priors can trigger optimization instability, hindering the formation of coherent reasoning pathways (DeepSeek-AI et al., 2025). (3) Outcome-Centric Filtering (Lv et al., 2025) relies on external proxies (e.g., accuracy) to estimate knowledge conflict. This creates an observational gap where external correctness masks latent shortcut learning, where the model attains answers via spurious cues rather than through faithful reasoning (Geirhos et al., 2020; Dziri et al., 2023). Consequently, these proxies fail to capture true model–data conflict, overlooking examples that require genuine structural adaptation (Dai et al., 2022). Across these regimes, the fundamental bottleneck is that data routing relies on coarse heuristics, such as pipeline order or outcome-based indicators, rather than intrinsic signals reflecting the model’s internal state.

Inspired by Schema Theory (Piaget, 1952), we posit that learning efficiency is fundamentally governed by the degree of conflict between new information and the model’s existing knowledge base. In this framework, compatible information is mastered through consolidation, while high-conflict information necessitates a fundamental restructuring of internal logic. To operationalize this principle, we adopt a gradient-based perspective, viewing gradients as the mathematical "feedback" signal derived from data. We propose that the distributional geometry of gradients serves as a critical proxy for this cognitive conflict. Specifically, we utilize statistical metrics such as the Gini coefficient to quantify gradient concentration. Prior studies suggest that concentrated updates (high Gini) typically correspond to data that deviates significantly from the model’s established knowledge base (Simsekli et al., 2019), whereas diffuse updates reflect global consistency with its current parametric state (Chizat et al., 2019). As illustrated in our ALFWorld case study (Figure 1), this gradient signal enables precise Data Arbitration. A "low-conflict" task (Case A) follows a standard routine and triggers diffuse gradients, making it an ideal candidate for SFT to consolidate behavioral norms. In contrast, a "high-conflict" scenario (Case B) involving complex error recovery generates highly concentrated gradients, signaling a failure of current logic that demands RL-driven exploration. By routing samples based on these intrinsic learning needs, **PRISM** (Partitioning Regimes via Internal Spatial-gradient Metrics) ensures both

training efficiency and robust generalization.

We evaluate PRISM on two challenging agent benchmarks: WebShop (online shopping) (Yao et al., 2022) and ALFWorld (embodied decision-making) (Shridhar et al., 2021). Empirical results demonstrate that PRISM achieves a Pareto improvement: it establishes a new state-of-the-art on ALFWorld (95.31) while reducing RL computational overhead by up to $3.22\times$. These results confirm that selective structural adaptation is both more robust and more efficient than exhaustive exploration. Notably, PRISM exhibits superior generalization capabilities across diverse backbones, including Qwen and Llama architectures. By precisely allocating high-conflict data to RL, the model achieves substantial performance gains in unseen environments, strongly validating the importance of distinguishing between consolidation and adaptation data for building robust agents.

Our contributions are summarized as follows:

- **Misalignment:** Formalizing the agent training bottleneck as functional SFT-RL mismatch from coarse-grained data allocation.
- **PRISM Framework:** A framework using spatial geometric structure of gradients to diagnose cognitive conflict, routing data between consolidation and adaptation regimes.
- **Efficiency:** We demonstrate SOTA performance across diverse backbones while delivering a $3.22\times$ training speedup via selective high-conflict allocation.

2 Methodology

We formalize PRISM, a framework designed to operationalize data routing by distinguishing between pattern consolidation and structural adaptation. We first establish the theoretical foundation by defining **gradient concentration** as a diagnostic proxy for **cognitive conflict**. Building upon this groundwork, we detail the implementation of PRISM via a three-stage pipeline (Fig. 2): (i) non-invasive gradient probing to capture the spatial geometric structure of learning dynamics; (ii) quantifying structural dissonance via statistical concentration metrics; and (iii) distribution-adaptive routing to allocate trajectories to their optimal learning regimes.

2.1 Gradient Concentration as a Proxy for Cognitive Conflict

This framework bridges the gap between low-level optimization dynamics and high-level cognitive

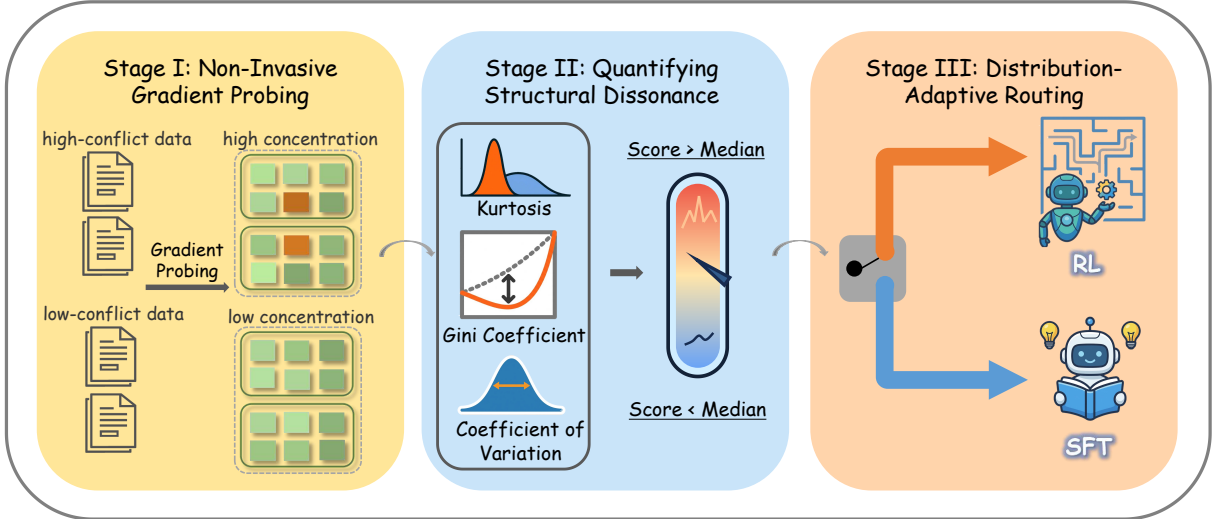


Figure 2: **Overview of PRISM.** The framework consists of three stages: (1) Non-Invasive Gradient Probing: Extracting update landscapes to capture internal reactions to each sample; (2) Quantifying Structural Dissonance: Measuring gradient concentration to diagnose the conflict between the data and existing knowledge; (3) Distribution-Adaptive Routing: Partitioning data based on concentration—samples with low-conflict (diffuse updates) are routed to SFT for consolidation, while those with high conflict (concentrated updates) are routed to RL for structural restructuring.

learning processes by treating the spatial geometric structure of gradients as a diagnostic signal for the conflict between input data and the model’s existing knowledge schema. This logic is rooted in the functional specificity of model parameters: research indicates that knowledge representation in Transformers is often localized within a sparse subset of specific “knowledge neurons” (Dai et al., 2022). When new information contradicts established patterns, the optimization process forces gradients into a significantly non-uniform spatial distribution, concentrating heavily within specific critical units to resolve internal logical dissonance (Simsekli et al., 2019; Meng et al., 2022). Consequently, this spatial concentration serves as an effective proxy for the magnitude of internal structural adaptation required by the model.

High Gradient Concentration signifies that the model must undergo intense localized updates to reconcile fundamental contradictions between the input and existing heuristics, signaling a regime of Structural Restructuring. In these high-conflict scenarios, Reinforcement Learning (RL) is indispensable, as its exploration mechanism drives the deep policy shifts necessary for logic realignment. Conversely, **Low Gradient Concentration** corresponds to diffuse, low-intensity parameter updates across the network, implying that new information can be seamlessly integrated without overhauling the underlying knowledge architecture (Chizat

et al., 2019). This state represents knowledge compatibility and pattern consolidation, where Supervised Fine-Tuning (SFT) provides an efficient and stable optimization path to refine behavioral norms at a lower computational cost and with minimal risk of destructive interference. By utilizing concentration metrics to distinguish between these internal regimes, PRISM enables principled **Data Arbitration**, selectively deploying RL only when structural adaptation is essential.

2.2 Stage I: Non-Invasive Gradient Probing

To capture the model’s internal reaction to new data without altering its weights, we perform Non-Invasive Gradient Probing. This stage serves as a lightweight “diagnostic scan” of the gradient landscape to identify the required optimization effort for each interaction trajectory τ_i .

Specifically, we utilize the ground-truth reference trajectories provided by the dataset as the supervisory signal. The rationale is to measure the dissonance between the model’s current policy and the expert behavior required by the task. We decompose the model’s parameter space Θ into N fine-grained functional units. For a Transformer with L layers, we focus on the linear projection matrices within the Attention and FFN blocks (e.g., $W_q, W_k, W_v, W_o, W_{\text{gate}}, \dots$), resulting in $N = 7L$ distinct parameter groups.

To eliminate confounds arising from varying se-

224 quence lengths, all input sequences are processed
 225 using a uniform, task-specific context length. We
 226 apply strict attention masking to ensure that gra-
 227 dients are computed solely on valid response to-
 228 kens, excluding padding artifacts. We execute a
 229 single backward pass without performing any op-
 230 timizer update, and then aggregate matrix-wise
 231 Frobenius norms to obtain a high-fidelity “snap-
 232 shot” of the internal learning state. This yields a
 233 high-dimensional gradient vector $\mathbf{g}_i \in \mathbb{R}_{\geq 0}^N$:

$$234 \quad \mathbf{g}_i = [\|\nabla_{\theta_1} \mathcal{L}(\tau_i)\|_F, \dots, \|\nabla_{\theta_N} \mathcal{L}(\tau_i)\|_F]^\top, \quad (1)$$

235 where $\mathcal{L}(\cdot)$ denotes the next-token prediction loss
 236 averaged over the valid response tokens, and
 237 $\|\nabla_{\theta_j} \mathcal{L}(\tau_i)\|_F$ is the Frobenius norm of the gradi-
 238 ent with respect to the j -th functional group θ_j .
 239 Significantly, this probing phase is computationally
 240 efficient and memory-friendly. By calculat-
 241 ing the norms on-the-fly during the backward pass
 242 and discarding the full gradient tensors, we main-
 243 tain a memory footprint nearly identical to a stan-
 244 dard forward pass, resulting in negligible overhead:
 245 the probing step accounts for only about 1–2% of
 246 the end-to-end wall-clock time of our full pipeline
 247 (probe+SFT+RL), as reported in Table 3.

248 2.3 Stage II: Quantifying Structural 249 Dissonance

250 While the raw gradient vector \mathbf{g}_i encodes both up-
 251 date intensity and structural shape, we prioritize
 252 the distributional shape to reveal the nature of the
 253 learning conflict, adhering to the principle of scale
 254 invariance. To robustly quantify this signal, we em-
 255 ploy a statistical concentration toolkit comprising
 256 three complementary metrics. The Gini Coeffi-
 257 cient measures the global inequality of gradient
 258 contributions across the network. Simultaneously,
 259 Kurtosis serves as a high-order diagnostic tool to
 260 detect heavy-tailed updates, identifying trajectories
 261 that force spiky adjustments in localized knowledge
 262 neurons while leaving the majority of functional
 263 circuits dormant. Finally, the Coefficient of Vari-
 264 ation (CV) captures relative instability by normal-
 265 izing dispersion against the mean update intensity.
 266 Together, these metrics triangulate cognitive dis-
 267 sonance from distinct statistical dimensions: high
 268 values signal an acute structural conflict requiring
 269 exploratory restructuring via RL, while low values
 270 indicate high data-model compatibility suitable for
 271 efficient consolidation via SFT.

272 In practice, for each trajectory τ_i , we compute a
 273 composite score $s_i = \phi(\mathbf{g}_i)$, where $\phi(\cdot)$ is a statis-

274 tical concentration operator (e.g., Gini) that maps
 275 the high-dimensional spatial geometric structure to
 276 a scalar value of cognitive dissonance.

277 2.4 Stage III: Distribution-Adaptive Routing

278 Finally, we partition the full corpus \mathcal{D} into disjoint
 279 subsets based on the quantified dissonance using
 280 a distribution-adaptive strategy. We compute the
 281 global statistics of the concentration scores \mathcal{S} and
 282 employ a non-parametric median split to define the
 283 routing boundary:

$$284 \quad \mathcal{D}_{\text{SFT}} = \{\tau_i \in \mathcal{D} \mid s_i \leq \text{Median}(\mathcal{S})\}, \quad (2)$$

$$285 \quad \mathcal{D}_{\text{RL}} = \{\tau_i \in \mathcal{D} \mid s_i > \text{Median}(\mathcal{S})\}, \quad (3)$$

286 where \mathcal{D} is the initial training corpus, $\mathcal{S} =$
 287 $\{s_1, \dots, s_{|\mathcal{D}|}\}$ represents the global set of compos-
 288 ite concentration scores for all trajectories, s_i
 289 is the structural dissonance score for the i -th trajec-
 290 tory, and $\text{Median}(\mathcal{S})$ serves as the threshold that
 291 dynamically partitions the corpus into consolida-
 292 tion (\mathcal{D}_{SFT}) and adaptation (\mathcal{D}_{RL}) regimes.

293 We employ a non-parametric median split as a
 294 robust, data-adaptive thresholding strategy. This
 295 approach ensures that data arbitration is grounded
 296 in the relative cognitive dissonance of the specific
 297 corpus, eliminating the need for per-task hyperpa-
 298 rameter tuning while maintaining a stable balance
 299 between plasticity and stability.

300 The rationale for selecting the median as the
 301 decision boundary is two-fold. First, this non-
 302 parametric boundary adaptively scales with the
 303 dataset’s inherent difficulty, ensuring that arbitra-
 304 tion is determined by the relative cognitive effort
 305 required by the model’s current state avoiding ar-
 306 bitrary constants. Second, it strikes a theoretical
 307 balance between stability and plasticity: routing
 308 too many samples to RL (a low threshold) intro-
 309 duces optimization noise from easy data, while
 310 routing too few (a high threshold) limits capacity
 311 for structural adaptation. This design choice is em-
 312 pirically validated in our Ablation Studies (Section
 313 4.3), where the median split consistently yields op-
 314 timal performance compared to extreme allocation
 315 ratios. Consequently, high-conflict trajectories are
 316 routed to RL for Structural Restructuring, while
 317 low-conflict trajectories are assigned to SFT for
 318 pattern consolidation, ensuring that each learning
 319 regime addresses the data’s intrinsic conflict.

Category	Method	Split (%)	Task-wise Success Rate (%)						
		SFT : RL	Pick	Look	Clean	Heat	Cool	Pick2	All
<i>(a) ALFWorld - Seen (In-Distribution)</i>									
Base	Base Model	- : -	39.29	7.14	0.00	0.00	0.00	0.00	9.38
Single Phase	SFT	100 : -	78.57	78.57	42.86	16.67	26.92	28.57	45.31
	GRPO	- : 100	80.65	28.57	86.21	52.94	66.67	57.89	67.19
	GiGPO	- : 100	87.10	42.86	75.86	76.47	61.11	52.63	69.53
Hybrid	Random	50 : 50	92.59	66.67	80.00	61.54	75.76	88.00	79.69
	HPT	50 : 50	92.59	77.78	92.59	75.00	75.00	75.00	85.16
	SFT-then-RL	100 : 100	96.67	72.73	96.15	85.71	80.00	86.36	88.28
PRISM (Ours)	Gini	50 : 50	92.68	100.00	100.00	100.00	95.00	91.67	95.31
	Kurtosis	50 : 50	95.12	88.89	100.00	100.00	75.00	87.50	91.41
	CV	50 : 50	95.12	88.89	100.00	90.00	85.00	75.00	89.84
<i>(b) ALFWorld - Unseen (Out-of-Distribution)</i>									
Base	Base Model	- : -	15.00	13.33	2.56	0.00	0.00	0.00	4.69
Single Phase	SFT	100 : -	70.00	20.00	15.38	4.17	10.00	0.00	19.53
	GRPO	- : 100	75.00	53.33	92.31	75.00	60.00	20.00	67.97
	GiGPO	- : 100	90.00	53.33	76.92	75.00	60.00	40.00	68.75
Hybrid	Random	50 : 50	45.00	33.33	74.36	54.17	90.00	65.00	60.94
	HPT	50 : 50	80.00	46.67	84.62	66.67	80.00	70.00	73.44
	SFT-then-RL	100 : 100	65.00	60.00	71.79	58.33	100.00	75.00	69.53
PRISM (Ours)	Gini	50 : 50	75.00	60.00	89.74	70.83	90.00	85.00	79.69
	Kurtosis	50 : 50	75.00	53.33	69.23	75.00	100.00	80.00	73.44
	CV	50 : 50	55.00	73.33	82.05	70.83	80.00	65.00	71.88

Table 1: Detailed Performance on ALFWorld (Qwen3-8B). Success Rate (%) across six task types on (a) Seen and (b) Unseen splits. SFT:RL indicates the allocation ratio. PRISM outperforms all hybrid baselines while requiring only half the data of the sequential SFT-then-RL pipeline.

3 Experiments

3.1 Experimental Setup

Benchmarks. Evaluation is conducted on two representative agentic benchmarks requiring distinct cognitive capabilities. **WebShop** (Yao et al., 2022), an interactive e-commerce environment, assesses the agent’s capacity for instruction following and attribute matching over long horizons, simulating real-world website navigation. Complementarily, **ALFWorld** (Shridhar et al., 2021) provides a text-based embodied simulation that demands compositional generalization for household tasks. Following standard protocols, performance is reported on both **Seen** (training distribution) and **Unseen** (generalization) splits to rigorously assess robustness against environmental shifts.

Implementation Details. Using **Qwen3-8B** (Yang et al., 2025) and **Llama-3.1-8B-Instruct** (AI@Meta, 2024) as backbones, we

implement SFT via **Llama-Factory** (Zheng et al., 2024) and RL (GRPO) via **verl-agent** (Feng et al., 2025). PRISM initiates with a **gradient probing phase** on the frozen base model to compute concentration metrics. These metrics serve as a filter to disentangle the dataset: low-conflict samples are assigned to SFT, while high-conflict samples are routed to RL. This selective allocation minimizes computational overhead while maximizing structural adaptation. See Appendix C for full hyperparameters.

Baselines We evaluate PRISM against three distinct categories: (1) **Monolithic Baselines** (100% budget): **SFT**, **GRPO**, and **GiGPO** (Feng et al., 2025), the current state-of-the-art method for agentic RL; (2) **Iso-Compute Baselines** (50/50 split): **Random** selection (serving as a control) and **HPT** (Lv et al., 2025), a leading outcome-aware method based on pass rates; and (3) **Canonical Pipeline: SFT-then-RL** (100%+100%), which

Method	Split (%)		Qwen3-8B		Llama-3.1-8B	
	SFT	RL	SR(%)	Score	SR(%)	Score
Base Model	-	-	17.19	28.78	1.56	1.56
SFT	100	-	42.97	73.87	18.75	29.34
GRPO	-	100	46.88	68.80	51.56	70.66
GiGPO	-	100	46.09	71.42	52.34	73.81
Random	50	50	55.47	78.21	54.69	79.40
HPT	50	50	54.69	75.48	55.47	80.46
SFT-then-RL	100	100	59.38	80.82	60.16	81.65
PRISM (Gini)	50	50	64.84	85.15	68.75	84.82
PRISM (Kurt)	50	50	63.28	83.87	64.06	81.79
PRISM (CV)	50	50	61.74	84.33	61.72	81.55

Table 2: Main Results on WebShop. We compare PRISM against baselines using different data allocation strategies. Data Split denotes the number of trajectories utilized for SFT and RL phases respectively. PRISM consistently outperforms the sequential baseline (SFT-then-RL) using only 50% of the total training budget, demonstrating the efficiency of concentration-aware data arbitration.

359 serves as a compute-intensive upper bound. Crucially, to ensure robustness, all reported results
360 represent the mean across three random seeds.
361

362 3.2 Main Results

363 Generalization and Efficiency on ALFWorld

364 As detailed in Table 1, PRISM demonstrates superior task mastery and generalization capabilities.
365 On **Seen** tasks (in-distribution), PRISM (Gini) achieves a remarkable success rate of **95.31%**,
366 significantly outperforming the sequential SFT-then-RL baseline (88.28%) and standard GRPO
367 (67.19%). Crucially, on **Unseen** tasks (out-of-distribution), PRISM exhibits exceptional robustness,
368 reaching a success rate of **79.69%**. This represents a substantial **10.16% absolute improvement**
369 over the sequential baseline.
370
371
372
373
374

375 These results validate our core hypothesis: routing low-conflict data to SFT facilitates efficient
376 pattern consolidation, while reserving high-conflict data for RL drives critical structural adaptation.
377 Unlike standard RL, which risks overfitting to environmental noise when trained on full data, PRISM
378 selectively targets trajectories requiring logical restructuring. Consequently, it achieves these gains
379 using only **50% of the RL compute budget**, effectively mitigating optimization interference and pre-
380 serving the model’s structural plasticity for novel scenarios.
381
382
383
384
385
386

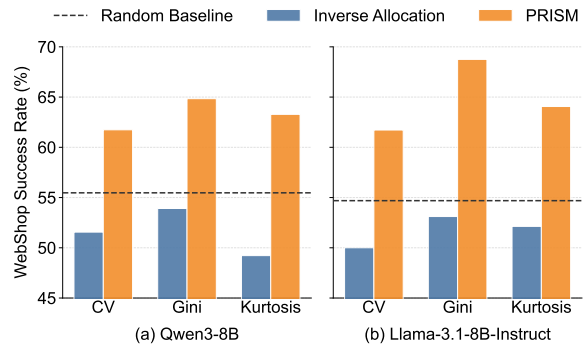


Figure 3: Ablation study of data routing strategies on WebShop. We compare PRISM (orange) with the Inverse allocation (blue: SFT on high-conflict data and RL on low-conflict data) under three concentration metrics for (a) Qwen3-8B and (b) Llama-3.1-8B-Instruct. The dashed line denotes the Random Baseline.

387 Backbone-Agnostic Robustness on WebShop

388 Table 2 highlights PRISM’s decisive advantage in
389 interactive decision-making across diverse model
390 architectures. On **Qwen3-8B**, PRISM (Gini) estab-
391 lishes a new state-of-the-art with a score of **85.15**
392 and a success rate of **64.84%**, surpassing both
393 the outcome-aware baseline HPT (75.48) and the
394 compute-intensive SFT-then-RL (80.82). Notably,
395 this superiority extends to **Llama-3.1-8B**, where
396 PRISM improves the Success Rate by **+8.59%** over
397 the sequential baseline (68.75% vs. 60.16%).

398 The consistent performance of the Gini metric
399 across both benchmarks suggests that **spatial gra-**
400 **di-ent concentration** serves as a robust proxy for
401 cognitive dissonance in agentic tasks. By filtering
402 out diffuse, template-like samples for SFT, PRISM
403 ensures that the expensive RL phase is exclusively
404 dedicated to resolving high-conflict bottlenecks
405 (e.g., complex attribute matching), thereby prevent-
406 ing the gradient dilution often observed in indis-
407 criminate full-data training.

408 4 Ablation and Analysis

409 4.1 Validating the Logic of Conflict-Aware 410 Routing

411 To validate the causal link, we partitioned data into
412 high- and low-concentration subsets using the Gini,
413 Kurtosis, and CV. We compared applying RL to
414 high-concentration and SFT to low-concentration
415 data against the reverse configuration (SFT on
416 high/RL on low). As shown in Figure 3, prioritizing
417 RL for high-concentration data significantly out-
418 performs both the random baseline and the reverse
419 setup. This confirms that concentrated updates sig-

Task	Method	Data Allocation	Training Wall-clock Time				Speedup
		(SFT : RL)	Probing	SFT Phase	RL Phase	Total Time	
WebShop	GRPO	0% : 100%	-	-	5h 53m	5h 53m	1.00×
	Random (50%)	50% : 50%	-	8m	3h 16m	3h 24m	1.73×
	PRISM (CV)	50% : 50%	1m 48s	8m	3h 09m	3h 18m	1.77×
	PRISM (Gini)	50% : 50%	1m 48s	8m	3h 10m	3h 19m	1.76×
	PRISM (Kurtosis)	50% : 50%	1m 48s	8m	2h 51m	3h 00m	1.95×
	GRPO	0% : 100%	-	-	36h 13m	36h 13m	1.00×
ALFWorld	Random (50%)	50% : 50%	-	7m	11h 41m	11h 48m	3.07×
	PRISM (CV)	50% : 50%	2m 16s	7m	11h 26m	11h 35m	3.12×
	PRISM (Gini)	50% : 50%	2m 16s	7m	11h 06m	11h 15m	3.22×
	PRISM (Kurtosis)	50% : 50%	2m 16s	7m	11h 11m	11h 20m	3.20×

Table 3: Computational Efficiency and Training Costs. Wall-clock time comparison on $8 \times$ NVIDIA A100 (80GB) GPUs. Data Allocation specifies the proportion of the dataset assigned to the SFT and RL phases, respectively. PRISM achieves superior results by intelligently partitioning a single dataset into optimal learning regimes, yielding a maximum speedup of **3.22**×

Routing Metric	WebShop (Score)	ALFWorld (SR %)
Gradient Magnitude (L_2)	79.75	90.63
PRISM (Spatial Gini)	85.15	95.31

Table 4: **Spatial Concentration vs. Gradient Magnitude.** We compare magnitude-based routing (allocating the top 50% of samples by magnitude to RL) against PRISM. Results show spatial concentration identifies structural adaptation requirements missed by gradient magnitude alone.

420 nal structural conflicts necessitating exploration,
421 whereas forcing RL on low-conflict data disrupts
422 consolidated norms (Chizat et al., 2019). Thus, se-
423 lective allocation based on concentration metrics is
424 empirically superior.

4.2 Disentangling Structural Conflict from Update Intensity

425 A critical question is whether PRISM simply prox-
426 ies sample difficulty. Comparing PRISM against
427 a Gradient Magnitude (routing top 50% samples
428 by L_2 norm to RL) in Table 4 reveals a decisive
429 advantage (+5.4% on WebShop). This distinction
430 is grounded in optimization dynamics: **High Mag-
431 nitude \neq RL Need.** Large gradient norms often
432 indicate simple “knowledge gaps” (e.g., unfamil-
433 iar entities) that are structurally compatible with
434 the model, making them ideal for efficient pattern
435 consolidation via SFT rather than expensive RL
436
437

438 exploration (Paul et al., 2021). Magnitude-based
439 routing inefficiently misallocates these learnable
440 samples to RL. In contrast, high concentration sig-
441 nals structural conflict. Concentrated updates im-
442 ply that the necessary logic correction is localized
443 within specific functional units (e.g., Knowledge
444 Neurons), reflecting a fundamental inconsistency
445 that requires the exploratory adaptation of RL to
446 resolve (Dai et al., 2022; Simsekli et al., 2019).

4.3 Sensitivity to Allocation Ratio

447 We evaluate PRISM’s sensitivity to routing thresh-
448 olds by varying the RL allocation ratio. As shown
449 in Figure 4, performance exhibits a distinct inverted
450 U-shape peaking near 50%. This revealing trend
451 highlights a critical trade-off: insufficient RL al-
452 location ($< 30\%$) provides inadequate structural
453 adaptation for high-conflict samples, while exces-
454 sive allocation ($> 70\%$) leads to gradient dilution.
455 Specifically, forcing RL on low-conflict data injects
456 exploratory noise into trivial behaviors, thereby
457 contaminating the gradients and interfering with
458 previously consolidated patterns.
459

4.4 Pareto Improvement: Optimization Disentanglement

460 Beyond raw performance, PRISM fundamentally
461 optimizes the computational utility of agent train-
462 ing. As detailed in Table 3, our framework reduces
463 wall-clock time by **nearly 50%** on WebShop and
464 achieves a **3.22**×

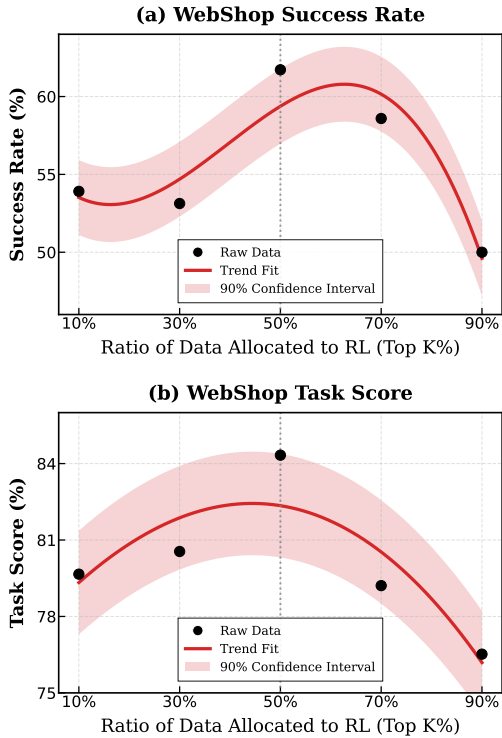


Figure 4: **Sensitivity to RL Allocation Ratio.** Performance of Qwen3-8B on WebShop (CV metric) across varying RL data proportions. The observed inverted U-shape peaks at a 50% split, indicating that a balanced allocation yields optimal performance compared to insufficient adaptation or excessive exploration.

results with the SFT-then-RL baseline reveals a critical insight regarding data scaling. While the sequential baseline processes 100% of the data via computationally expensive RL to achieve competitive results, PRISM achieves superior performance using only $\sim 50\%$ of the RL budget. This indicates that blindly forcing exploration on well-consolidated knowledge yields diminishing returns.

PRISM effectively disentangles the optimization objectives: it delegates pattern consolidation to the cost-efficient SFT phase, while reserving the expensive RL budget for trajectories requiring structural restructuring of key neurons. This synergy realizes a Pareto improvement in the performance-efficiency trade-off, proving that smarter data arbitration is superior to exhaustive exploration.

5 Related Works

Data Allocation in SFT-RL. Data allocation balances imitation and exploration, yet current paradigms often rely on coarse heuristics that ignore the interplay between model state and data difficulty. Monolithic sequencing (Ouyang et al., 2022) uses rigid schedules, failing to distinguish

between pattern consolidation and structural adaptation (Zhou et al., 2023). Similarly, universal exploration (Shao et al., 2024; Feng et al., 2025) indiscriminately applies RL, which can trigger optimization instability on high-conflict data lacking SFT-consolidated priors (DeepSeek-AI et al., 2025). While outcome-centric filtering (Lv et al., 2025) uses external proxies, it suffers from an “observational gap” where correctness masks latent shortcuts or unfaithful reasoning (Geirhos et al., 2020; Dziri et al., 2023). In contrast, PRISM shifts to internal learning dynamics, utilizing the spatial geometric structure of gradients to quantify intrinsic conflict. This enables granular arbitration based on actual cognitive demand rather than rigid pipeline orders or outcome-based heuristics.

Gradient-Based Diagnostics Gradients provide high-fidelity diagnostics for internal dynamics and functional specialization (Zhao et al., 2025; Dai et al., 2022). Unlike intensity- or similarity-based approaches (Paul et al., 2021; Zhao et al., 2024), PRISM leverages the spatial geometric structure of updates, aligning with mechanistic evidence of localized representations (Geva et al., 2021). Specifically, concentrated updates signal structural restructuring to resolve logical inconsistencies (Simsekli et al., 2019; Meng et al., 2022), while diffuse updates reflect knowledge compatibility and consolidation (Chizat et al., 2019; Agarwal et al., 2022). We repurpose these signals into a proactive arbitration mechanism for optimal regime routing.

6 Conclusion

In this work, we introduced **PRISM**, a framework that bridges cognitive learning principles with neural optimization to resolve the long-standing data arbitration challenge in agent training. By utilizing the spatial geometric structure of gradients as an intrinsic diagnostic for cognitive conflict, PRISM effectively disentangles the training process into pattern consolidation via SFT and structural adaptation via RL. Our results confirm that precision in data routing outweighs raw volume: PRISM not only establishes new state-of-the-art benchmarks but also mitigates optimization interference, yielding superior generalization. This approach represents a significant **Pareto improvement**, achieving these gains with a $3.22\times$ training speedup. Ultimately, PRISM marks a shift from heuristic-based pipelines toward a principled, dynamics-aware orchestration of LLM post-training.

7 Limitations

Despite its robust performance and efficiency, PRISM has several limitations that warrant further exploration. First, due to **computational constraints**, our evaluation is primarily focused on 7B–8B parameter models. While we hypothesize that the **spatial geometric structure** of gradients is a scale-invariant mechanistic property of Transformers, extensive verification on large-scale models (e.g., 70B+ parameters) remains for future work. Second, we currently employ a **static routing strategy** based on initial gradient concentration to isolate diagnostic signals and minimize computational overhead. This approach does not account for the dynamic evolution of a model’s internal state, where a high-conflict sample might transition into a routine consolidation candidate as training progresses. Finally, our scope is concentrated on **agentic decision-making** benchmarks. While these tasks effectively highlight the functional divergence between SFT and RL, the generalizability of our gradient-based diagnostic to other complex domains, such as advanced mathematical reasoning or open-ended creative generation, requires further empirical investigation.

References

Chirag Agarwal, Daniel D’souza, and Sara Hooker. 2022. Estimating example difficulty using variance of gradients. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10368–10378.

AI@Meta. 2024. [Llama 3 model card](#).

Lenaic Chizat, Edouard Oyallon, and Francis Bach. 2019. On lazy training in differentiable programming. *Advances in neural information processing systems*, 32.

Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. 2025. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493–8502.

DeepSeek-AI and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Nouha Dziri and 1 others. 2023. Faith and fate: Limits of transformers on compositionality. *arXiv preprint arXiv:2305.18654*.

Lang Feng, Zhenghai Xue, Tingcong Liu, and Bo An. 2025. Group-in-group policy optimization for llm agent training. *arXiv preprint arXiv:2505.10978*.

Robert Geirhos and 1 others. 2020. Shortcut learning in deep neural networks. *Nature Machine Intelligence*.

Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Xingtai Lv, Yuxin Zuo, Youbang Sun, Hongyi Liu, Yuntian Wei, Zhekai Chen, Lixuan He, Xuekai Zhu, Kaiyan Zhang, Bingning Wang, and 1 others. 2025. Towards a unified view of large language model post-training. *arXiv preprint arXiv:2509.04419*.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. *Advances in neural information processing systems*, 35:17359–17372.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.

Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. 2021. Deep learning on a data diet: Finding important examples early in training. *Advances in neural information processing systems*, 34:20596–20607.

Jean Piaget. 1952. *The origins of intelligence in children*, volume 8. International Universities Press New York.

Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiushi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. 2025. Toolrl: Reward is all tool learning needs. *arXiv preprint arXiv:2504.13958*.

Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, and 1 others. 2024. Toolllm: Facilitating large language models to master 16000+ real-world apis. In *ICLR*.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

647	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	A Details of Concentration Metrics	702
648	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	In this section, we provide the formal definitions for	703
649	Zhang, YK Li, Yang Wu, and 1 others. 2024.	the gradient concentration metrics used to quantify	704
650	Deepseekmath: Pushing the limits of mathematical	cognitive dissonance.	705
651	reasoning in open language models. <i>arXiv preprint</i>		
652	<i>arXiv:2402.03300</i> .		
653	Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté,	Gradient Vector Construction. For a given tra-	706
654	Yonatan Bisk, Adam Trischler, and Matthew	jectory τ_i , let $\mathcal{L}(\tau_i)$ denote the standard next-token	707
655	Hausknecht. 2021. Alfworld: Aligning text and	prediction loss, averaged over all valid tokens in	708
656	embodied environments for interactive learning. In	the sequence. To characterize the spatial geo-	709
657	<i>ICLR</i> .	metric structure of the model’s internal response,	710
658	Umut Simsekli, Levent Sagun, and Mert Gurbuzbala-	we analyze the gradients with respect to the spe-	711
659	ban. 2019. A tail-index analysis of stochastic gradi-	cific linear projection weights of the Transformer	712
660	ent noise in deep neural networks. In <i>International</i>	backbone. For a model with L layers, we de-	713
661	<i>Conference on Machine Learning</i> , pages 5827–5837.	fine the parameter groups for the l -th layer as	714
662	PMLR.	$\mathcal{P}_l = \{W_q, W_k, W_v, W_o, W_{\text{gate}}, W_{\text{up}}, W_{\text{down}}\}$. Ag-	715
663	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	gregating across all layers, we obtain a total of	716
664	Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,	$N = 7L$ parameter groups. This multi-layered	717
665	Chengen Huang, Chenxu Lv, Chujie Zheng, Dayi-	grouping allows us to capture the distribution of	718
666	heng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge,	optimization effort across the network’s functional	719
667	Haoran Wei, Huan Lin, Jialong Tang, and 41 oth-	units, providing the necessary resolution to mea-	720
668	ers. 2025. Qwen3 technical report. <i>arXiv preprint</i>	sure spatial concentration.	721
669	<i>arXiv:2505.09388</i> .	We define the gradient concentration vector $\mathbf{g}_i \in$	722
670	Shunyu Yao, Howard Chen, John Yang, and Karthik	$\mathbf{R}_{\geq 0}^N$ as the collection of Frobenius norms for each	723
671	Narasimhan. 2022. Webshop: Towards scalable real-	parameter group’s gradient matrix:	724
672	world web interaction with grounded language agents.		
673	<i>Advances in Neural Information Processing Systems</i> ,	$\mathbf{g}_i = [\ \nabla_{\theta_1} \mathcal{L}(\tau_i)\ _F, \dots, \ \nabla_{\theta_N} \mathcal{L}(\tau_i)\ _F]^\top. \quad (4)$	725
674	35:20744–20757.	Let μ_i and σ_i denote the arithmetic mean and stan-	726
675	Minxing Zhang, Yi Yang, Roy Xie, Bhuwan Dhingra,	dard deviation of the elements in \mathbf{g}_i , respectively.	727
676	Shuyan Zhou, and Jian Pei. 2025. Generalizability	ϵ is a small constant ($1e^{-8}$) added for numerical	728
677	of large language model-based agents: A comprehen-	stability.	729
678	sive survey. <i>arXiv preprint arXiv:2509.16330</i> .		
679	Yang Zhao, Li Du, Xiao Ding, Yangou Ouyang, Hep-	1. Gini Coefficient. The Gini coefficient mea-	730
680	eng Wang, Kai Xiong, Jinglong Gao, Zhouhao Sun,	sures the inequality of the gradient contribution	731
681	Dongliang Xu, Qing Yang, and 1 others. 2025. Be-	distribution. We first sort the elements of \mathbf{g}_i in	732
682	yond similarity: A gradient-based graph method for	non-decreasing order , such that $g_{i,(1)} \leq g_{i,(2)} \leq$	733
683	instruction tuning data selection. In <i>Proceedings</i>	$\dots \leq g_{i,(N)}$. The metric is calculated as:	734
684	<i>of the 63rd Annual Meeting of the Association for</i>		
685	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	$s_i^{\text{Gini}} = \frac{\sum_{j=1}^N (2j - N - 1) g_{i,(j)}}{N \sum_{j=1}^N g_{i,(j)} + \epsilon}. \quad (5)$	735
686	pages 24391–24404.	A higher Gini coefficient indicates that a small	736
687	Yang Zhao, Li Du, Xiao Ding, Kai Xiong, Zhouhao Sun,	subset of parameter groups dominates the gradient	737
688	Shi Jun, Ting Liu, and Bing Qin. 2024. Decipher-	updates (sparse activation), suggesting structural	738
689	ing the impact of pretraining data on large language	conflict.	739
690	models through machine unlearning. In <i>Findings of</i>		
691	<i>the Association for Computational Linguistics: ACL</i>	2. Kurtosis. We employ the Fourth Standard-	740
692	<i>2024</i> , pages 9386–9406.	ized Moment (Pearson’s Kurtosis) to quantify the	741
693	Yaowei Zheng, Richong Zhang, Junhao Zhang, and 1	“tailedness” of the gradient distribution. This serves	742
694	others. 2024. Llamafactory: Unified efficient fine-	as a detector for extreme outliers in optimization	743
695	tuning of 100+ language models. <i>arXiv preprint</i>	pressure. Given the large number of parameter	744
696	<i>arXiv:2403.13372</i> .		
697	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer,		
698	Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping		
699	Yu, Lili Yu, and 1 others. 2023. Lima: Less is more		
700	for alignment. <i>Advances in Neural Information Pro-</i>		
701	<i>cessing Systems</i> , 36:55006–55021.		

groups ($N \gg 100$), we utilize the population formula without small-sample bias correction:

$$s_i^{\text{Kurt}} = \frac{1}{N} \sum_{j=1}^N \left(\frac{g_{i,j} - \mu_i}{\sigma_i + \epsilon} \right)^4 - 3. \quad (6)$$

High kurtosis implies that the gradients are characterized by infrequent but extreme updates, distinguishing "spiky" structural adaptation signals from Gaussian noise.

3. Coefficient of Variation (CV). The Coefficient of Variation provides a normalized measure of concentration, describing the extent of variability in relation to the mean of the population:

$$s_i^{\text{CV}} = \frac{\sigma_i}{\mu_i + \epsilon}. \quad (7)$$

This metric captures the relative instability of the update signal, serving as a robust proxy for global model dissonance.

B Qualitative Analysis of Routed Trajectories

To validate the cognitive dissonance hypothesis, we manually inspected trajectories routed to distinct phases.

- **SFT-Routed (Low Concentration):** Typically involve straightforward instruction following or keyword matching (e.g., "Click the 'Search' button"). The model's priors are sufficient, resulting in diffuse gradients.
- **RL-Routed (High Concentration):** Involve counter-intuitive reasoning or correcting a previous error (e.g., ALFWorld: "The apple is not in the fridge, checking the cabinet"). These induce concentrated updates as specific attention heads must be re-weighted to shift the search strategy.

C Implementation Details

Gradient Probing Configuration To ensure consistency between the diagnostic and training phases, the Non-Invasive Gradient Probing (Stage I) utilizes the same context length constraints as the subsequent RL stage. Specifically, input sequences are standardized to a fixed length of 2048 tokens for ALFWorld and 4096 tokens for WebShop. Sequences exceeding these limits are truncated, while shorter ones are padded with strict masking applied during gradient computation to avoid padding bias.

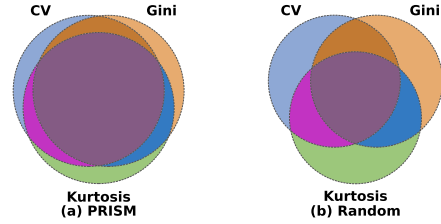


Figure 5: **Venn Diagram of Data Selection Consensus.** The intersection shows that approximately **60%** of the high-conflict trajectories are consistently identified by all three statistical metrics. This high degree of convergence significantly exceeds the **12.5%–25.0%** expected from random overlapping splits, demonstrating that PRISM captures a stable underlying structural dissonance signal regardless of the specific concentration metric employed.

SFT We implement the SFT stage using the LLaMA-Factory framework. We perform full-parameter fine-tuning on Qwen3-8B for 3 epochs using the AdamW optimizer. The learning rate is initialized at 1×10^{-5} with a cosine decay schedule and a warmup ratio of 0.1. We employ a per-device batch size of 4 with 4 gradient accumulation steps, training in bfloat16 precision.

RL For our method, we employ the GRPO algorithm without KL divergence penalties and set the rollout size to 8. We adopt the environment configurations and reward structures from the GiGPO framework. Specifically, the actor learning rate is set to 1×10^{-6} . We use a rule-based reward function: +10 for success, 0 for failure, and a penalty of -0.1 for invalid actions. Consistent with the probing phase, we limit prompts to **2048 tokens** for ALFWorld and **4096 tokens** for WebShop, with a maximum of 50 environment steps per episode for ALFWorld and 15 for WebShop. For the GiGPO baseline reported in our experiments, we strictly follow the original hyperparameter settings provided in (Feng et al., 2025).

D Consensus Analysis of Gradient Concentration Metrics

As illustrated in Figure 5, we observe a substantial overlap among the high-conflict subsets identified by these metrics. This empirical evidence suggests that while individual metrics may align more closely with specific task dynamics, they largely converge on a core set of high-conflict data. This consensus indicates that PRISM captures a universal and robust signal of cognitive dissonance, rather

than being an artifact of specific metric selection.

E Experimental Environments and Task Decomposition

We evaluate our framework on two complex agent benchmarks: **WebShop** and **ALFWorld**. These environments require the agent to demonstrate diverse capabilities, ranging from navigating e-commerce interfaces to solving interactive household tasks.

E.1 WebShop

WebShop simulates an e-commerce website environment, requiring models to navigate interfaces to find and purchase products that match specific user attributes.

Evaluation Metrics. Following the standard protocol of the WebShop benchmark (Yao et al., 2022), we evaluate the performance of our agent using two primary metrics: **Average Score** and **Success Rate (SR)**.

- **Average Score:** This metric measures the granularity of task completion by calculating the attribute overlap between the product purchased by the agent and the user’s instruction. For each episode i , the environment computes a scalar score $S_i \in [0, 1]$, which is a weighted sum of rewards based on four dimensions: product category matching, attribute recall, option selection accuracy, and price constraints. Formally, the score is calculated as:

$$S_i = \text{TypeScore} \times \left(\frac{N_{\text{attr}} + N_{\text{option}} + \mathbb{I}_{\text{price}}}{N_{\text{total}}} \right) \quad (8)$$

where N_{attr} and N_{option} denote the number of matched attributes and options respectively, and $\mathbb{I}_{\text{price}}$ is an indicator function for price satisfaction. We report the mean score averaged across all test episodes.

- **Success Rate (SR):** This is a stricter metric evaluating the agent’s ability to perfectly satisfy user goals. An episode is considered successful if and only if the agent achieves a perfect score (i.e., $S_i = 1.0$). This implies that the purchased item meets all specified criteria, including correct category, attributes, options, and price limits. SR denotes the percentage of episodes where the agent successfully completed the task.

E.2 ALFWorld

ALFWorld aligns TextWorld with the ALFRED benchmark, consisting of interactive household tasks that require multi-step reasoning and decision-making.

Task Decomposition. We report results across six ALFWorld sub-task categories: Pick (single-object retrieval), Look (object search/navigation), Clean (cleaning appliances), Heat (heating state transitions), Cool (cooling state transitions), and Pick2 (two-object pick-and-place).

Evaluation Metrics. Similar to WebShop, we report the **Success Rate (SR)** for ALFWorld. An episode is considered successful if the agent completes the goal state within the maximum number of steps. We report both the overall SR and the task-wise SR for the six categories mentioned above.

F Theoretical Motivation: Why High-Conflict Trajectories Benefit from RL Exploration

The main text argues that high gradient concentration indicates a structural mismatch between the current policy and the target behavior, motivating routing such trajectories to RL. This appendix provides a mechanistic explanation for why exploration-based, group-relative RL (e.g., GRPO) is well-matched to this regime. We present the argument as an intuition consistent with policy-gradient learning dynamics, rather than as a formal equivalence between gradient concentration under SFT and the RL training signal.

1. High conflict tends to create distinct rollout modes under sampling. When a state-action decision is aligned with the model’s current behavior, stochastic sampling from π_θ often produces similar trajectories with small qualitative variation. In contrast, under structural mismatch, the policy is more likely to admit multiple competing action modes for the same state (e.g., relying on superficial heuristics versus executing a faithful reasoning chain). As a result, sampling can expose distinct outcome patterns (success/failure, or different intermediate behaviors), creating the diversity necessary for trial-and-error refinement in policy optimization (Schulman et al., 2017).

2. GRPO is most informative when there is within-group contrast. GRPO-style learning constructs its update from relative comparisons

914 within a sampled group of trajectories (e.g., group-
915 relative advantages), rather than from matching a
916 single reference trace (Shao et al., 2024; Feng et al.,
917 2025). This implies a simple requirement: the sam-
918 pled group must contain meaningfully different
919 outcomes for the relative signal to be discrimina-
920 tive.

- 921 • **Low-conflict regime: limited contrast**
922 **yields weakly discriminative relative feed-**
923 **back.** For consolidated trajectories, sampled
924 rollouts tend to be homogeneous in outcomes
925 and rewards. In this case, group-relative nor-
926 malization/ranking provides little separation
927 between trajectories, so the relative learning
928 signal becomes less informative and can be
929 sensitive to stochasticity without yielding sys-
930 tematic improvement (Shao et al., 2024; Feng
931 et al., 2025).
- 932 • **High-conflict regime: outcome separation**
933 **enables contrastive credit assignment.** Un-
934 der structural mismatch, sampling is more
935 likely to produce both better and worse roll-
936 outs with distinct reward profiles. This within-
937 group separation makes group-relative up-
938 dates informative: the optimizer can assign
939 credit by reinforcing behaviors that lead to
940 verified success and suppressing those leading
941 to failure, without requiring imitation of a sin-
942 gle fixed trace (Shao et al., 2024; Feng et al.,
943 2025).

944 **3. Exploration supports selective policy shifts**
945 **where imitation can be brittle.** A key advantage
946 of routing high-conflict trajectories to RL is that
947 exploration allows the learner to search over alter-
948 native behaviors and update the policy selectively
949 based on feedback, rather than forcing the model
950 to reproduce a particular trajectory. This is con-
951 sistent with observations that RL post-training can
952 induce behavioral improvements beyond SFT-only
953 pipelines by leveraging reward-driven feedback to
954 shape policy updates (Ouyang et al., 2022; Guo
955 et al., 2025). In PRISM, this motivates concentrat-
956 ing RL budget on trajectories that exhibit structural
957 mismatch, while using SFT to consolidate already-
958 compatible behaviors.

959 **G Robustness Analysis: Architecture** 960 **Invariance and Confound Control**

961 A potential concern in gradient-based analysis is
962 whether the varying sizes of parameter matrices

(e.g., W_{down} vs. W_q) introduce confounds in the
963 concentration metrics. We argue that PRISM is
964 robust to these variations due to **Architecture In-**
965 **variance.** 966

967 While larger parameter matrices naturally yield
968 larger gradient norms, this introduces a **constant**
969 **systematic bias** rather than a data-dependent vari-
970 able. Since the model architecture remains static,
971 this bias affects all trajectories identically. More-
972 over, our chosen metrics (e.g., Gini Coefficient)
973 are theoretically **scale-invariant**—multiplying a
974 subset of dimensions by a constant factor preserves
975 the relative inequality score, effectively canceling
976 out layer-wise scaling artifacts.

977 To empirically validate this, we conducted a
978 **sensitivity analysis** by normalizing the gradient
979 norms by the square root of the parameter count
980 ($\|\mathbf{g}\|_F / \sqrt{N_{param}}$). We observed that this normal-
981 ization resulted in **nearly identical data rankings**
982 (Spearman’s $\rho > 0.99$ on both benchmarks) com-
983 pared to the raw Frobenius norms. This confirms
984 that PRISM’s median-split routing is driven by gen-
985 uine structural conflict rather than architectural di-
986 mensions.

987 **H AI Assistance Disclosure**

988 We acknowledge the use of AI tools solely for lan-
989 guage polishing and grammatical editing to im-
990 prove the readability of this manuscript. All sci-
991 entific claims, experimental data, and empirical
992 results were rigorously verified by the human au-
993 thors to ensure authenticity and accuracy.