

REPOSHAPLEY: Shapley-Enhanced Context Filtering for Repository-Level Code Completion

Anonymous ACL submission

Abstract

Repository-level code completion benefits from retrieval-augmented generation (RAG). However, controlling cross-file evidence is difficult because chunk utility is often interaction-dependent: some snippets help only when paired with complementary context, while others harm decoding when they conflict. We propose REPOSHAPLEY, a coalition-aware context filtering framework supervised by Shapley-style marginal contributions. Our module **ChunkShapley** constructs offline labels by (i) single-chunk probing with teacher-forced likelihood to estimate signed, weighted effects, (ii) a surrogate game that captures saturation and interference, (iii) exact Shapley computation for small retrieval sets, and (iv) bounded post-verification that selects a decoding-optimal coalition using the frozen generator. We distill verified *KEEP* or *DROP* decisions and retrieval triggering into a single model via discrete control tokens. Experiments across benchmarks and backbones show that REPOSHAPLEY improves completion quality while reducing harmful context and unnecessary retrieval. Code: <https://anonymous.4open.science/r/a7f3c9>.

1 Introduction

Repository-level code completion requires resolving non-local dependencies across files, including project-specific APIs, shared contracts, and invariants (Jimenez et al., 2024; Ding et al., 2024b). Retrieval-Augmented Generation (RAG) addresses this setting by injecting cross-file evidence into Code LMs (Lewis et al., 2020; Kang et al., 2024; Shrivastava et al., 2023; Bairi et al., 2023). However, effective retrieval control remains a bottleneck. Under a fixed context budget, the model must identify truly useful evidence from a noisy candidate pool, where many chunks are redundant and some are actively misleading (Ding et al., 2024a; Zhang et al., 2023; Wei et al., 2025; Liu et al.,

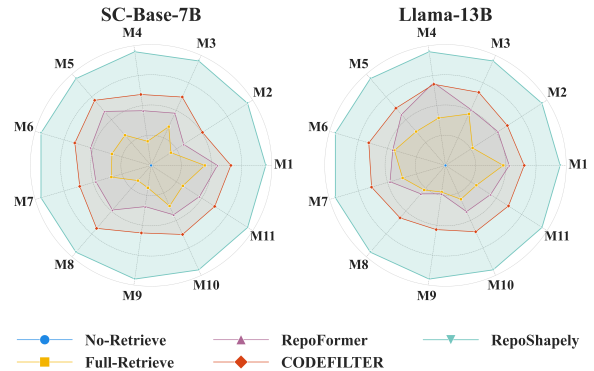


Figure 1: Performance radar charts on StarCoder-Base-7B and CodeLlama-13B. The plots display relative improvements over the No-Retrieve baseline (center). REPOSHAPLEY achieves SOTA performance across 11 tested metrics (As shown in table 1).

2024a; Yoran et al.).

The core difficulty is that chunk utility is often interaction-dependent. A snippet may appear uninformative in isolation yet become decisive when paired with complementary context, such as an interface declaration together with its implementation. Conversely, a plausible chunk can degrade generation when it co-occurs with conflicting evidence, such as deprecated versus updated APIs (Shi et al., 2023; Xu et al., 2024). Therefore, methods that score candidates independently can misestimate the utility of the multi-chunk context that is actually consumed at test time (Khandelwal et al., 2020; Yan et al., 2024; Bertsch et al., 2025).

To address this, we adopted a coalition-first approach. Retrieval control should be supervised by signals that reflect how a chunk behaves within a set, rather than in isolation. We introduce REPOSHAPLEY, a framework that learns to filter context using Shapley-style marginal contributions.

Our approach has two stages. First, we propose **ChunkShapley**, an offline labeling pipeline for interaction-aware supervision. Considering that computing Shapley values directly with the genera-

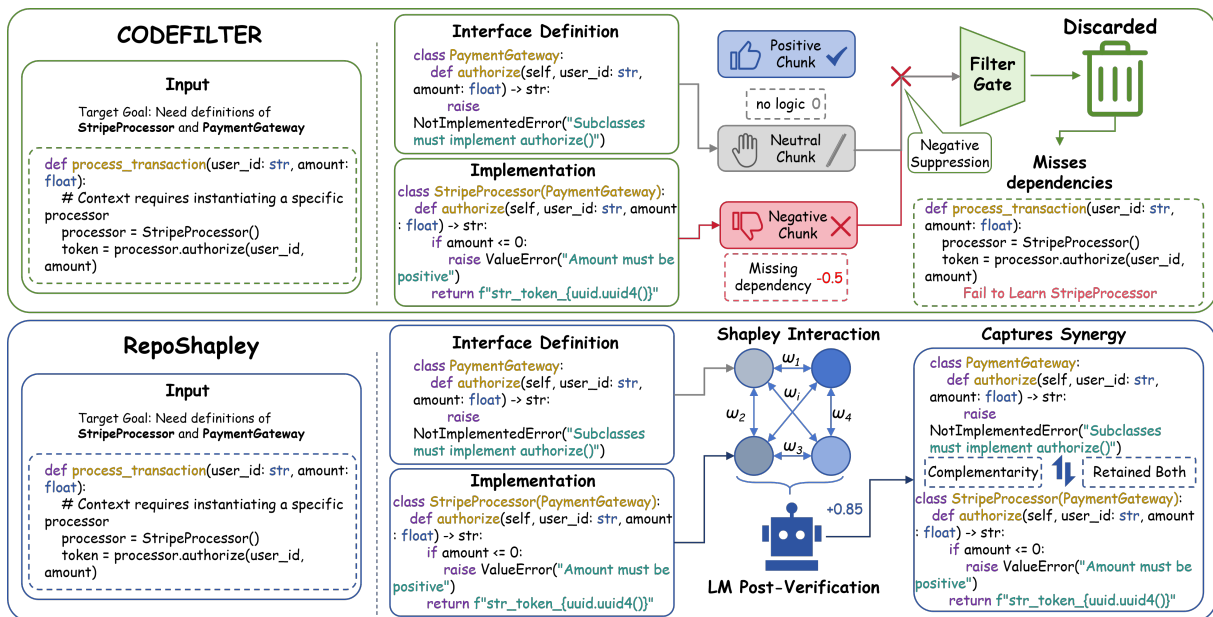


Figure 2: Under the same input context and the exact same retrieved candidate chunks, CODEFILTER makes decisions from independent per-chunk signals and can break under interaction effects, whereas REPOSHAPLEY performs coalition-aware filtering that more reliably removes high-score noise while preserving complementary evidence.

tor is prohibitive, we introduce a structured logistic surrogate that can capture saturation and conflict efficiently. We then apply a verification step to ground the selected coalitions in the generator’s actual decoding behavior. Second, we distill the resulting coalition-derived labels into a single generator via discrete control tokens, which we call REPOSHAPLEY. This distillation enables efficient, interaction-aware retrieval control at inference time. As shown in Figure 1, REPOSHAPLEY achieves SOTA performance across benchmarks and metrics, supporting our motivation that coalition-aware supervision is crucial for difficult cross-file completion. Our contributions are as follows:

- **Coalition-aware supervision for context filtering.** We formulate context selection as a cooperative game and use Shapley marginal contributions to capture complementarity and conflict beyond independent scoring.
- **ChunkShapley: Practical shapley labeling for chunk filtering.** We combine single-chunk probing with a structured surrogate utility to compute exact Shapley values on small retrieval sets ($K=10$). We further select a verified coalition from a bounded candidate pool under decoding-time metrics.
- **REPOSHAPLEY: distillation for online retrieval control.** We distill verified keep and

drop decisions into discrete control tokens, enabling a single model to decide when to retrieve and which chunks to keep.

2 Related Work

Repository-Level RAG and Retrieval Control. RAG mitigates non-local dependencies in code completion by retrieving cross-file evidence (Lewis et al., 2020; Izacard and Grave, 2021; Parvez et al., 2021; Guu et al., 2020; Jiang et al., 2023; Mallen et al., 2023; Yao and Fujita, 2024). Recent work improves context quality via iterative retrieval (Gao et al., 2023; Zhang et al., 2023; Shrivastava et al., 2023; Zhang et al., 2025), structure-aware indexing, including dataflow or call graphs (Cheng et al., 2024; Liu et al., 2024c), and dedicated benchmarks (Ding et al., 2023; Liu et al., 2024b; Li et al., 2025; Generation, 2024; Yang et al., 2025). In parallel, retrieval control has received increasing attention, focusing on when to retrieve and what to retain under a fixed context budget. RepoFormer (Wu et al., 2024) triggers retrieval through self-evaluation, while CODEFILTER (Li et al., 2025) filters chunks using independent likelihood-based signals. However, these controllers largely assess chunks in isolation. As a result, they do not explicitly account for combinatorial interactions such as complementarity between interfaces and implementation. In contrast, we cast context filtering as a coalition scoring problem to model such inter-

dependencies.

Shapley Values in RAG and Supervision. Shapley values (Shapley et al., 1953) provide an axiomatic notion of marginal contribution and have been widely used in interpretability, including SHAP-style formulations (Lundberg and Lee, 2017; Ghorbani and Zou, 2019; Sundararajan et al., 2017). In RAG, prior work applies Shapley-style analysis to attribute outputs to retrieved documents (Nematov et al., 2025; Ye and Yoganarasimhan, 2025) or to estimate token-level importance (Asai et al., 2024; Xiao et al., 2025). Our use differs in its role in the pipeline. Instead of post-hoc analysis of a frozen system, we use Shapley-style marginalization to construct supervision for retrieval control. We then distill the resulting coalition reasoning into a token-level policy, enabling practical retrieval decisions during generation.

3 Methodology

3.1 Repository-level Retrieval-Augmented Code Completion

Repository-level code completion requires grounding generation in cross-file information such as project-specific APIs, shared utilities, and type or contract conventions. RAG addresses this by retrieving candidate snippets from the repository. However, retrieved evidence is often interaction-heavy: a snippet may be useful only when paired with complementary context, and seemingly relevant snippets can degrade generation when they introduce conflicting implementations.

Problem setup. Given a repository \mathcal{R} and a target file, each instance is represented as $(X_{\text{in}}, X_{\text{out}}, Y)$. Here $X_{\text{in}} = (X_p, X_s)$ is the in-file context in fill-in-the-middle (FIM) format with prefix X_p and suffix X_s , X_{out} denotes a cross-file pool constructed from other files in \mathcal{R} , and Y is the ground-truth missing span between X_p and X_s (Zhang et al., 2023; Wu et al., 2024).

Retrieval and generation. A retriever R queries X_{out} with X_{in} and returns top- K candidate chunks $X_{\text{cc}} = R(X_{\text{in}}, X_{\text{out}}) = \{cc_1, \dots, cc_K\}$. A generator G_θ then predicts the completion \hat{Y} conditioned on X_{in} and a selected subset $X_S \subseteq X_{\text{cc}}$. Hence, the key problem is to estimate chunk utility and retain the subset that best supports generating Y .

3.2 Interaction-aware Chunk Attribution via Shapley Values

Why independent chunk scoring is insufficient. Retrieved code snippets rarely contribute independently. A chunk can be uninformative on its own but become essential when paired with complementary context such as an interface and its implementation. Conversely, a seemingly relevant snippet may reduce generation quality when it conflicts with other retrieved evidence. As a result, per-chunk scores computed in isolation can be a poor proxy for the utility of the multi-chunk context used at test time.

Subset utility as a cooperative game. We therefore evaluate chunks at the set level. Given top- K candidates, we treat each chunk as a player and any subset as a coalition. Let $D = 1, \dots, K$ index candidates and $S \subseteq D$ denote a coalition, with $X_S = cc_i : i \in S$. We define the coalition value as the normalized teacher-forced log-likelihood gain on the ground-truth completion:

$$v(S | X_{\text{in}}, Y) = \ell(X_{\text{in}}, X_S) - \ell(X_{\text{in}})$$

$$\ell(C) = \frac{1}{|Y|} \log p_\theta(Y | C).$$

where $\log p_\theta(Y | C) = \sum_{t=1}^{|Y|} \log p_\theta(y_t | y_{<t}, C)$. By construction, $v(\emptyset | X_{\text{in}}, Y) = 0$, and $v(S)$ can be negative when retrieved context decreases model likelihood. Appendix C.4 compares log-likelihood with metric-based utilities (EM/ES) and shows log-likelihood yields the best downstream performance.

Shapley attribution. We quantify interaction-aware chunk contributions using the Shapley value (Shapley et al., 1953), which is defined as the average marginal gain of chunk i over all coalitions:

$$\phi_i = \sum_{S \subseteq D \setminus \{i\}} \frac{|S|! (K - |S| - 1)!}{K!} \Delta v_i(S)$$

$$\Delta v_i(S) = v(S \cup \{i\} | X_{\text{in}}, Y) - v(S | X_{\text{in}}, Y).$$

Intuitively, $\phi_i > 0$ indicates that chunk i is helpful on average across different co-occurring contexts, while $\phi_i \leq 0$ suggests redundancy or harm under interactions. Shapley values satisfy *efficiency*: $\sum_{i \in D} \phi_i = v(D | X_{\text{in}}, Y)$, allowing negative attributions when some chunks reduce coalition utility.

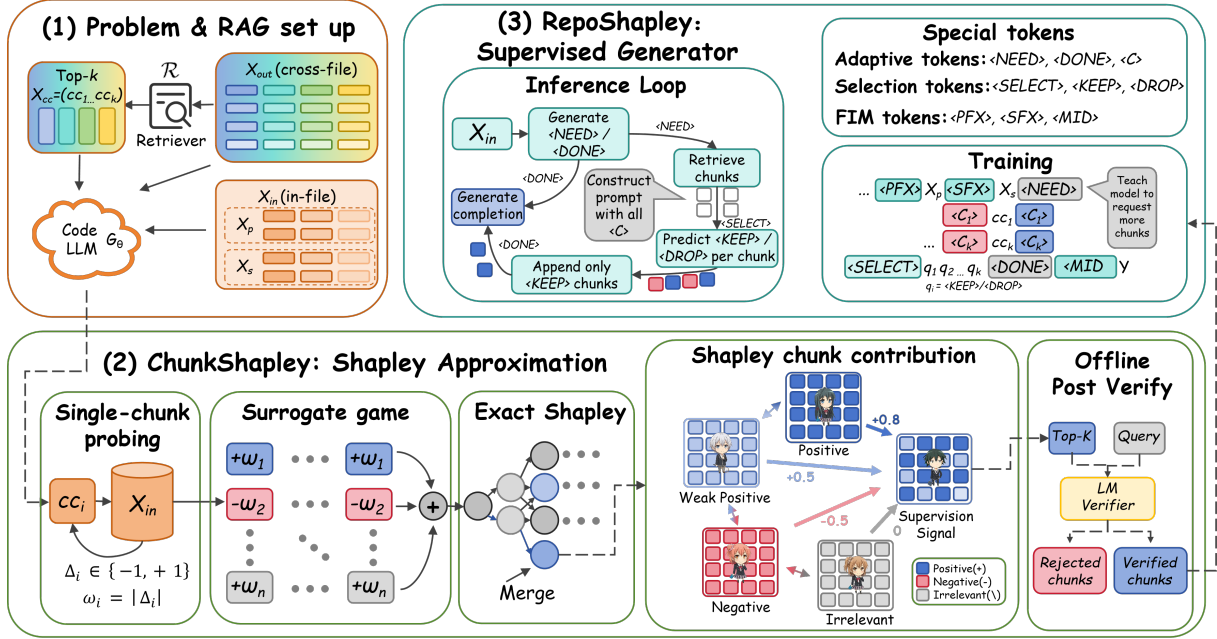


Figure 3: The overall framework of REPOSHAPLEY. The pipeline consists of two phases: (2) An offline ChunkShapley module that estimates the interaction-aware contribution of each chunk; and (3) An online Shapley-supervised Generator LLM trained to control retrieval and filter contexts based on the estimated Shapley values.

3.3 ChunkShapley: Practical Shapley Labeling for Chunk Filtering

Exact Shapley computation under the true coalition utility $v(\cdot)$ is impractical, as it would require evaluating the generator on exponentially many subsets. We therefore propose **ChunkShapley**, an *offline* labeling pipeline that (a) probes each chunk once to obtain a signed effect, (b) defines a lightweight surrogate game to approximate interaction patterns, (c) computes exact Shapley values under the surrogate by enumerating all 2^K coalitions, which is inexpensive since $v_{\text{sur}}(\cdot)$ is closed-form, and (d) performs bounded post-verification with the frozen generator to ground final keep and drop labels in decoding-time behavior. Algorithmic details are deferred to Appendix Alg. 2.

(a) Single-chunk probing. We first compute a per-instance baseline score using teacher forcing and probe each candidate in isolation. Let $\ell(C)$ denote the normalized teacher-forced log-likelihood. For each retrieved chunk cc_i , we define its single-chunk effect

$$\begin{aligned} \Delta_i &= \ell(X_{in}, \{cc_i\}) - \ell(X_{in}) \\ y_i &= \text{sign}(\Delta_i), \quad \omega_i = |\Delta_i|. \end{aligned}$$

To ensure consistent likelihood estimation under a limited context window, we preserve the full tar-

get span Y and apply **left-truncation** only to the input context (i.e., X_{in} and retrieved chunks).

(b) Logistic surrogate game. While ranking by Δ_i captures individual relevance, it ignores coalition dynamics. To model interactions efficiently, we define a one-dimensional surrogate utility. Given (y_i, ω_i) , we aggregate coalition S via a weighted vote:

$$g(S) = \sum_{i \in S} \omega_i y_i, \quad v_{\text{sur}}(S) = \sigma(\beta g(S)) - \sigma(0). \quad 247$$

where $\sigma(\cdot)$ is the sigmoid and $\beta > 0$ controls the saturation scale. This surrogate is not meant to match the full combinatorial utility; it targets two dominant effects for filtering. The sigmoid yields diminishing returns: when $|g(S)|$ is large, $\sigma'(\beta g(S)) \approx 0$, so additional similarly-signed evidence contributes little, capturing redundancy under a fixed budget. Conflicts are expressed by negative votes ($y_i = -1$), which reduce $g(S)$ and can suppress $v_{\text{sur}}(S)$ even when some chunks are individually helpful. By the way, subtracting $\sigma(0)$ ensures $v_{\text{sur}}(\emptyset) = 0$ and keeps utilities centered. The surrogate remains lightweight for exhaustive subset evaluation, while any residual mismatch to decoding-time behavior is addressed by verification.

264 **(c) Exact Shapley values under the surrogate.**

265 We compute Shapley values using the subset form
266 under the surrogate utility:

$$267 \phi_i = \frac{1}{K} \sum_{S \subseteq D \setminus \{i\}} \frac{v_{\text{sur}}(S \cup \{i\}) - v_{\text{sur}}(S)}{\binom{K-1}{|S|}}.$$

268 Since our $v_{\text{sur}}(S)$ is closed-form, evaluating all 2^K
269 subsets is computationally negligible for small re-
270 trieval sizes ($K \leq 10$). This allows us to obtain
271 *exact* Shapley values under v_{sur} , avoiding the vari-
272 ance of sampling approximations.

273 In contrast, computing interactions using the
274 heavy generator G_θ would require exponentially
275 many coalition evaluations and is intractable.
276 Therefore, we use ϕ_i under the surrogate as a pro-
277 posal signal and rely on post-verification to finalize
278 the decision.

279 **(d) Post-verification via a bounded candi-
280 date pool.**

281 Because decoding quality is non-
282 monotonic in context, so the positive attributions
283 alone do not guarantee improved greedy decoding.
284 Since the surrogate is only a proxy, we verify a
285 small candidate pool with the frozen generator and
286 select the coalition that maximizes decoding-time
287 quality. This step is used only for offline label con-
288 struction with access to Y ; the inference never uses
289 Y .

289 Let π_ϕ and π_Δ be indices sorted by ϕ_i and Δ_i .
290 We build a de-duplicated set \mathcal{C} containing: (i) Shap-
291 ley prefixes $\{\pi_\phi[1:n]\}_{n=1}^{N_v}$, (ii) short Δ prefixes as
292 a strong single-chunk baseline, and (iii) size-2/3
293 combinations among top- L chunks by Δ to explic-
294 itly probe local synergies. For each $S \in \mathcal{C}$, we
295 decode with the frozen generator and choose

$$296 S^* = \arg \max_{S \in \mathcal{C}} (\text{ES}(\hat{Y}_S, Y), \text{EM}(\hat{Y}_S, Y))$$

297 using lexicographic maximization (ES first, EM as
298 tie-break). We then treat S^* as the teacher keep/-
299 drop labels for distillation.

300 **Verified labels for retrieval triggering.** The
301 post-verification step also yields an oracle deci-
302 sion on whether retrieval is necessary. Let \hat{Y}_\emptyset
303 be the decoding result using only in-file context
304 X_{in} , and let \hat{Y}_{S^*} be the decoding result using the
305 verification-selected coalition S^* . We define the
306 retrieval-control label as

$$307 r^* = \begin{cases} \langle \text{DONE} \rangle, & \text{if } \text{ES}(\hat{Y}_{S^*}, Y) - \text{ES}(\hat{Y}_\emptyset, Y) \leq \epsilon \\ \langle \text{NEED} \rangle, & \text{otherwise.} \end{cases}$$

308 where ϵ is a small margin tuned on the validation
309 set (default $\epsilon = 0$ unless stated otherwise). This
310 label is used only for offline supervision; inference
311 never accesses Y .

312 **3.4 REPOSHAPLEY: Distilling ChunkShapley
313 into Signal Tokens**

314 While ChunkShapley provides robust coalition-
315 aware supervision, the pipeline is too computationally
316 intensive for online use. We therefore propose
317 REPOSHAPLEY, which *distills* verified coalition
318 decisions into discrete control tokens, enabling a
319 single generator to efficiently decide *when* to re-
320 trieve and *which* chunks to retain at inference time.

321 **Signal tokens and verified labels.** We introduce
322 retrieval-control tokens $\mathcal{T}_R = \{\langle \text{NEED} \rangle, \langle \text{DONE} \rangle\}$ to
323 decide whether cross-file evidence is required, and
324 candidate-selection tokens $\mathcal{T}_S = \{\langle \text{KEEP} \rangle, \langle \text{DROP} \rangle\}$
325 to indicate which retrieved chunks should be re-
326 tained.

327 Step (d) outputs a *verification-selected coalition*
328 S^* by evaluating a small set of Shapley-proposed
329 candidate coalitions \mathcal{C} using the frozen gener-
330 ator under decoding-time constraints, to match
331 decoding-time behavior. We treat S^* as the teacher
332 keep/drop label set and distill it into token-level
333 supervision by assigning, for each retrieved chunk
334 cc_i ,

$$335 Q(cc_i) = \begin{cases} \langle \text{KEEP} \rangle & \text{if } i \in S^* \\ \langle \text{DROP} \rangle & \text{otherwise.} \end{cases}$$

336 In this way, surrogate Shapley signals are used only
337 to propose promising coalitions, while the student
338 model learns to imitate the *verified* coalition-level
339 behavior encoded by S^* , turning combinatorial sub-
340 set selection into a single-shot, controllable genera-
341 tion policy at inference time.

342 **Training: two-format verbalized supervision.**

343 Following the standard separation of evidence
344 selection and completion generation in retrieval-
345 augmented code modeling, we train a single model
346 with two serialized views of each instance. Format-
347 1 supervises *selection*: given the in-file context and
348 retrieved candidates, the model emits a keep/drop
349 decision token for each chunk. Format-2 super-
350 vises *generation*: the model produces the missing
351 span conditioned only on the kept evidence. Both
352 formats reuse the same control tokens and share
353 all parameters, enabling the model to learn selec-
354 tion and generation within a unified autoregressive
355 interface.

Format-1: Selection. Given the in-file context and the retrieved candidate list, the model predicts a length- K decision sequence $q_{1:K} \in \{\langle \text{KEEP} \rangle, \langle \text{DROP} \rangle\}^K$ under a dedicated $\langle \text{SELECT} \rangle$ marker. Let $[X_p]$ and $[X_s]$ denote tokenized FIM prefix and suffix, and let $\text{Pack}(X_{cc})$ be the deterministic serialization of retrieved candidates $X_{cc} = \{cc_1, \dots, cc_K\}$:

$$\text{Pack}(X_{cc}) = \langle C_1 \rangle [cc_1] \langle /C_1 \rangle \cdots \langle C_K \rangle [cc_K] \langle /C_K \rangle.$$

The Format-1 sequence is

$$\text{F1} : \langle \text{PFX} \rangle [X_p] \langle \text{SFX} \rangle [X_s] \langle \text{NEED} \rangle \\ \text{Pack}(X_{cc}) \langle \text{SELECT} \rangle q_1 q_2 \cdots q_K \langle \text{DONE} \rangle.$$

We supervise q_i using the *verified teacher coalition* S^* : $q_i^* = \langle \text{KEEP} \rangle$ if $i \in S^*$ and $\langle \text{DROP} \rangle$ otherwise.

Format-2: Generation. To teach the model how to complete code *given filtered evidence*, we construct a generation format that includes only the chunks in S^* and then decodes the target span in FIM mode:

$$\text{F2} : \langle \text{PFX} \rangle [X_p] \langle \text{SFX} \rangle [X_s] \langle \text{NEED} \rangle \text{Pack}(C_{S^*}) \\ \langle \text{DONE} \rangle \langle \text{MID} \rangle [Y].$$

No-retrieval format. If retrieval is unnecessary ($r^* = \langle \text{DONE} \rangle$), we drop the cross-file block and the selection head:

$$\langle \text{PFX} \rangle [X_p] \langle \text{SFX} \rangle [X_s] \langle \text{DONE} \rangle \langle \text{MID} \rangle [Y].$$

This indicates that in-file context suffices.

Remark. We reuse $\langle \text{DONE} \rangle$ both as the retrieval decision token and as a block delimiter; the two usages are unambiguous from their fixed positions in the sequence. The retrieval-control token ($\langle \text{NEED} \rangle / \langle \text{DONE} \rangle$) is learned with teacher forcing as a next-token target (counted in \mathcal{L}_R), rather than provided as an oracle input.

Objectives with masked contexts. We mask all in-file and cross-file *content* tokens in the loss and compute gradients only on *generated targets* (control tokens, selection tokens, and the completion Y). Let $r^* \in \mathcal{T}_R = \{\langle \text{NEED} \rangle, \langle \text{DONE} \rangle\}$ be the retrieval-control label. For retrieval-needed instances (Formats F1/F2), $r^* = \langle \text{NEED} \rangle$; for no-retrieval instances, $r^* = \langle \text{DONE} \rangle$. For Format-1, we optimize retrieval triggering and selection:

$$\mathcal{L}_R^{\text{F1}} = -\log P_G(r^* | X_{\text{in}}) \\ \mathcal{L}_S^{\text{F1}} = -\sum_{i \in \mathcal{J}} \log P_G(q_i^* | X_{\text{in}}, X_{\text{cc}}, r^*, q_{<i}^*) \\ \mathcal{L}^{\text{F1}} = \lambda_R \mathcal{L}_R^{\text{F1}} + \lambda_S \mathcal{L}_S^{\text{F1}}$$

Algorithm 1: REPOSHAPLEY Inference

Process

Input: Generator G , Retriever R ,
Cross-file pool X_{out} , In-file context
 $X_{\text{in}} = (X_p, X_s)$;

Token sets $\mathcal{T}_R = \{\langle \text{NEED} \rangle, \langle \text{DONE} \rangle\}$,
 $\mathcal{T}_S = \{\langle \text{KEEP} \rangle, \langle \text{DROP} \rangle\}$; threshold t_c .

Output: Completed code \hat{Y} .

```

1  $X \leftarrow (\langle \text{PFX} \rangle, X_p, \langle \text{SFX} \rangle, X_s)$ 
2  $r \leftarrow \text{Select}(\text{Softmax}_{\mathcal{T}_R}(G(\cdot | X)), t_c)$ 
3 if  $r = \langle \text{DONE} \rangle$  then
4    $X \leftarrow \text{append}(X, \langle \text{MID} \rangle)$ 
5   return  $\hat{Y} \leftarrow G(X)$ 
6 end
7  $X_{\text{cc}} \leftarrow R(X_{\text{in}}, X_{\text{out}})$ 
8  $X_{\text{sel}} \leftarrow X \oplus \langle \text{NEED} \rangle \oplus \text{Pack}(X_{\text{cc}}) \oplus \langle \text{SELECT} \rangle$ 
9  $(q_1, \dots, q_K) \leftarrow G(X_{\text{sel}})$ 
10  $\hat{S} \leftarrow \{i \in \{1, \dots, K\} : q_i = \langle \text{KEEP} \rangle\}$ 
11  $X \leftarrow X \oplus \langle \text{NEED} \rangle \oplus \text{Pack}(C_{\hat{S}}) \oplus \langle \text{DONE} \rangle$ 
12  $X \leftarrow \text{append}(X, \langle \text{MID} \rangle)$ 
13 return  $\hat{Y} \leftarrow G(X)$ 

```

where $\mathcal{J} = \{1, \dots, K\}$ for retrieval-needed instances and $\mathcal{J} = \emptyset$ for no-retrieval instances.

For Format-2, we optimize retrieval triggering and generation conditioned on the verified filtered context. Let $X_{S^*} = \text{Pack}(C_{S^*})$ denote the serialized filtered evidence.

$$\mathcal{L}_R^{\text{F2}} = -\log P_G(r^* | X_{\text{in}}) \\ \mathcal{L}_Y^{\text{F2}} = -\sum_{t=1}^T \log P_G(y_t | y_{<t}, X_{\text{in}}, X_{S^*}, r^*) \\ \mathcal{L}^{\text{F2}} = \lambda_R \mathcal{L}_R^{\text{F2}} + \mathcal{L}_Y^{\text{F2}}.$$

Here \mathcal{L}_R is implemented as the cross-entropy on the next-token prediction at the retrieval-control position (i.e., immediately after $\langle \text{SFX} \rangle [X_s]$), rather than a separate classifier.

During training, we either (i) include both formats for each instance, or (ii) sample one format per instance with a fixed mixing ratio. The final objective is the expectation over the chosen format:

$$\mathcal{L} = \mathbb{E}_{F \sim \pi} [\mathcal{L}^F], \quad F \in \{\text{F1}, \text{F2}\}.$$

Inference. At inference time, REPOSHAPLEY makes retrieval decisions in one autoregressive rollout. Given the in-file context, the model first predicts a retrieval-control token $r \in \mathcal{T}_R = \{\langle \text{NEED} \rangle, \langle \text{DONE} \rangle\}$. If $r = \langle \text{DONE} \rangle$, it directly performs FIM decoding to generate the completion.

Table 1: Code completion performance in the Infilling setting.

Model	Strategy	RepoEval						CCLongEval			CCEval	
		Line		API		Function		Chunk		Func	Line	
		EM (M1)	ES (M2)	EM (M3)	ES (M4)	UT (M5)	ES (M6)	EM (M7)	ES (M8)	ES (M9)	EM (M10)	ES (M11)
SC-Base-1B	No-Retrieve	43.14	67.39	38.03	66.81	21.67	47.29	30.62	60.54	47.16	18.72	42.85
	Full-Retrieve	52.27	73.13	44.18	69.09	25.61	55.93	37.49	64.04	50.72	22.38	47.26
	RepoFormer	54.71	76.52	45.73	72.41	28.46	57.69	41.93	70.21	54.37	25.42	49.18
	CODEFILTER	57.19	78.84	48.37	75.66	31.13	59.91	44.52	72.48	56.59	27.81	52.03
	REPOSHAPLEY	61.34 +4.15	82.78 +3.94	53.62 +5.25	79.53 +3.87	35.84 +4.71	64.39 +4.48	48.57 +4.05	77.52 +5.04	61.18 +4.59	32.26 +4.45	56.37 +4.34
SC-Base-3B	No-Retrieve	48.12	72.38	40.17	68.91	24.93	51.52	36.16	65.19	49.63	21.82	45.58
	Full-Retrieve	57.84	77.21	48.83	72.68	30.58	58.16	42.61	68.29	53.84	25.92	50.31
	RepoFormer	58.59	79.16	49.82	74.63	32.89	60.62	46.38	72.11	56.39	28.85	52.16
	CODEFILTER	61.21	81.09	51.97	77.62	35.18	63.26	49.62	74.58	58.51	30.84	55.29
	REPOSHAPLEY	64.93 +3.72	85.27 +4.18	56.38 +4.41	81.72 +4.10	39.91 +4.73	68.16 +4.90	53.52 +3.90	78.83 +4.25	62.84 +4.33	35.79 +4.95	59.41 +4.12
SC-Base-7B	No-Retrieve	51.62	75.51	43.83	71.29	25.62	52.71	38.91	66.62	52.84	23.37	48.01
	Full-Retrieve	58.26	77.79	50.38	75.01	32.26	60.21	44.62	69.19	55.16	28.51	52.91
	RepoFormer	59.83	79.26	51.31	77.46	35.71	61.19	46.84	74.16	57.11	29.62	55.49
	CODEFILTER	61.49	81.41	53.62	79.29	37.79	63.41	49.16	77.26	59.84	32.11	57.84
	REPOSHAPLEY	65.81 +4.32	86.59 +5.18	58.79 +5.17	84.11 +4.82	41.84 +4.05	68.16 +4.75	54.73 +5.57	81.29 +4.03	64.62 +4.78	36.59 +4.48	62.91 +5.07
Llama-7B	No-Retrieve	51.89	73.42	41.53	66.98	24.81	44.56	37.21	65.16	50.37	18.16	43.34
	Full-Retrieve	60.18	78.91	48.76	73.16	29.93	52.21	45.41	69.37	52.11	23.41	47.46
	RepoFormer	60.52	79.36	49.31	75.91	33.19	52.64	46.84	69.56	52.16	24.26	48.31
	CODEFILTER	63.76	82.31	52.62	78.54	32.84	54.49	50.76	74.91	54.74	27.16	51.68
	REPOSHAPLEY	68.31 +4.55	86.76 +4.45	57.28 +4.66	83.11 +4.57	37.56 +4.72	59.14 +4.65	55.21 +4.45	79.19 +4.28	59.41 +4.67	31.23 +4.07	55.24 +3.56
Llama-13B	No-Retrieve	53.81	74.84	42.19	67.96	26.31	47.14	41.91	67.46	52.71	20.97	45.88
	Full-Retrieve	61.41	79.29	49.81	77.41	31.69	54.21	47.36	70.61	55.24	25.53	50.20
	RepoFormer	62.19	81.51	50.46	79.21	34.39	54.44	48.96	71.11	55.41	27.20	52.20
	CODEFILTER	64.16	82.71	53.01	78.99	35.41	57.76	51.31	74.19	58.81	29.91	54.69
	REPOSHAPLEY	68.89 +4.73	87.11 +4.40	57.66 +4.65	83.41 +4.42	40.11 +4.70	62.39 +4.63	55.91 +4.60	78.59 +4.40	63.51 +4.70	35.05 +5.14	59.37 +4.68

If $r = \langle \text{NEED} \rangle$, we retrieve K cross-file candidates $X_{cc} = cc_1, \dots, cc_K$ and serialize them as $\text{Pack}(X_{cc})$. Conditioned on this packed block, the model outputs a length- K selection sequence under $\langle \text{SELECT} \rangle$, where $(q_1, \dots, q_K) \in \mathcal{T}_S^K$ and $\mathcal{T}_S = \langle \text{KEEP} \rangle, \langle \text{DROP} \rangle$. We then keep only chunks with $q_i = \langle \text{KEEP} \rangle$, append them to the prompt, and generate \hat{Y} via FIM decoding after emitting $\langle \text{MID} \rangle$. Alg. 1 provides the full procedure. We use \oplus to denote token sequence concatenation.

4 Experiments

4.1 Experimental Setup

Dataset. We curate 290k Python repositories from The Stack (Kocetkov et al., 2023) after strict quality filtering (LOC constraints, AST parsing, and deduplication; Appendix A.2). Following (Wu et al., 2024), we sample 7.5k repositories to construct 50k labeled training instances: for each instance, we retrieve top-10 cross-file chunks using Jaccard similarity (Jaccard, 1912) and assign supervision derived from ChunkShapley. During data labeling, we discard instances whose verification-selected coalition S^* fails to reach a minimum completion quality, so that $\text{ES}(\hat{Y}_{S^*}, Y) < \tau_{\text{es}}$, to ensure supervision reliability. We split the data into 95%/5% for training and validation.

Models and Training. We fine-tune StarCoderBase (SCB-1B/3B/7B) (Li et al., 2023) and CodeLlama (Llama-7B/13B) (Roziere et al., 2023) for 2

epochs using a learning rate of 2×10^{-5} with linear decay and 5% warm-up. We set $\lambda_R = \lambda_S = 2.0$, max sequence length to 4096. With a global batch size of 512 on 8 NVIDIA H100 (80GB), training takes on average 2.2/6.5/15.4 hours for SCB-1B/3B/7B and 15.8/28.6 hours for Llama-7B/13B, respectively. (Details are shown in Appendix. B)

Benchmarks and Metrics. We evaluate on three repository-level code completion benchmarks: RepoEval (Zhang et al., 2023), CrossCodeEval (Ding et al., 2023), and CrossCodeLongEval (Wu et al., 2024). Together they cover line, API, chunk, and function-level completion tasks under realistic cross-file dependencies. We consider two prompting settings: **Infilling** (FIM with $X_{\text{in}} = (X_p, X_s)$) and **Left-to-right** (prefix-only with $X_{\text{in}} = X_p$). Following prior work (Wu et al., 2024), we report Exact Match (EM) and Edit Similarity (ES) for non-function tasks, and unit-test pass rate (UT) for function tasks (Formulations are shown in Appendix. A.1).

Baselines. We compare REPOSHAPLEY against: (1) **No-Retrieve** (in-file only); (2) **Full-Retrieve** (Zhang et al., 2023) (top-10 sparse retrieval); (3) **RepoFormer** (Wu et al., 2024) (selective retrieval); and (4) **CODEFILTER** (Li et al., 2025) (likelihood-based filtering). CODEFILTER serves as the primary baseline to highlight the benefit of interaction-aware supervision.

Table 2: Component Ablation of REPOSHAPLEY on RepoEval. We investigate components in (A) Labeling and (B) Distillation. Baseline is SC-Base-1B.

Method / Variant	RepoEval-Line		RepoEval-API		Latency
	EM	ES	EM	ES	(ms/req)
RepoFormer	54.71	79.26	45.73	72.41	661
CODEFILTER	57.19	78.84	48.37	75.66	947
REPOSHAPLEY	61.34	82.78	53.62	79.53	1053
A. Labeling Strategy					
1. w/o Post-verification	38.50	54.44	36.15	55.81	–
2. Δ -only labeling	58.45	77.12	48.46	75.26	–
3. Linear Surrogate	59.92	76.41	50.73	77.09	–
4. Uniform Weights	60.18	80.97	51.82	77.38	–
B. Distillation					
5. Format-1 only	5.56	8.27	2.34	5.66	523
6. Format-2 only	59.88	79.11	52.12	77.49	830
7. No Trigger	61.26	81.33	52.15	78.81	1462

4.2 Main Results

Tables 1 and 3 show that REPOSHAPLEY consistently improves repository-level infilling across benchmarks and backbones, validating our core hypothesis that supervision derived from evidence coalitions better reflects interaction-heavy retrieval.

First, interaction-blind filtering remains brittle. While adaptive controllers generally outperform Full-Retrieve, methods trained from per-chunk labels (CODEFILTER) can still overfit to isolated similarity and fail to account for complementarity and conflict that only appear when multiple chunks are concatenated. This gap is most visible on harder settings that require resolving non-local dependencies such as Function, where selecting the right combination of evidence matters more.

Second, coalition-aware supervision yields the strongest gains on difficult tasks. On SC-Base-7B, REPOSHAPLEY improves RepoEval API from 53.62/79.29 to **58.79/84.11** (EM/ES) and raises RepoEval Function unit-test pass rate from 37.79 to **41.84**, outperforming CODEFILTER by clear margins. These improvements align with our motivation: modeling evidence interactions helps retain complementary context while suppressing conflicting or redundant chunks.

Finally, the gains generalize beyond RepoEval. REPOSHAPLEY also delivers consistent improvements on long-context and chunk-level benchmarks, on CCLongEval Chunk it improves ES from 77.26 to **81.29** on SC-Base-7B and from 74.19 to **78.59** on Llama-13B, indicating that the learned keep and drop policy transfers across evaluation granularities and context regimes.

Although REPOSHAPLEY introduces additional

computation during offline labeling, its inference-time overhead remains modest. As shown in Table 2, REPOSHAPLEY runs at 1053 ms/req, which is comparable to CODEFILTER (947 ms) and within the same runtime scale as RepoFormer (661 ms).

4.3 Ablation Study & Analysis

We study how each component of REPOSHAPLEY affects performance by ablating (A) the offline labeling pipeline and (B) the online distillation strategy on StarCoderBase-1B with RepoEval (Table 2).

Coalition-aware labeling matters. Ablations in Part A show that modeling interactions is necessary for reliable filtering. Using Shapley signs without post-verification (Row 1) causes a large drop, indicating that signed marginal effects alone are not stable under prerequisite dependencies. Replacing coalition-based attribution with single-chunk probing (Row 2) also hurts performance, suggesting that independent scores miss synergy among chunks. Simplifying the surrogate utility by removing the sigmoid (Row 3) or using uniform weights (Row 4) further degrades results, supporting our design for capturing saturation and conflict effects.

Distillation and triggering improve inference. Part B shows that the training design is essential. Training with selection-only signals (Row 5) fails to produce usable code, while generation-only training (Row 6) lags behind the full model due to residual noise. Removing the trigger (Row 7) yields similar accuracy but increases latency, confirming that the learned trigger reduces unnecessary retrieval while maintaining generation quality.

5 Conclusion

In this work, we study repository-level retrieval control under strong evidence interaction effects. We propose ChunkShapley, a practical approximation that combines K single-chunk probes with a continuous logistic surrogate game to compute Shapley-style attributions efficiently, and further introduces a compact post-verification step that grounds selection in the true generator’s behavior. This design bridges the gap between interaction-aware attribution and deployable retrieval control: the surrogate Shapley stage provides a principled ranking that accounts for complementary or conflicting evidence, while post-verification mitigates non-monotonicity and metric-specific failure modes by selecting the best coalition from a small, structured candidate pool.

561 Limitations

562 Our method has several limitations. First, the surrogate utility is constructed from single-chunk probes, which provides a tractable approximation but can miss higher-order interactions where multiple individually weak chunks become useful only jointly. Our bounded post-verification mitigates this issue, yet it cannot guarantee recovering such cases outside the candidate pool. Second, our offline labeling requires enumerating all 2^K subsets under the surrogate game; although we use small retrieval sizes ($K=10$), the cost grows exponentially and may limit scalability to larger candidate sets or higher budgets. Third, the verification stage is tied to a particular decoding setup and evaluation criteria (greedy decoding with ES/EM); the selected coalition may vary under different decoding algorithms, stochastic sampling, or task-specific objectives. Finally, the labeling pipeline involves multiple teacher-forced forward passes for probing and additional decoding runs for verification, which increases offline computation and can be costly when constructing large-scale supervision.

584 References

585 Anthropic. 2025. [System card: Claude opus 4 & claude sonnet 4](#). Technical report, Anthropic. Accessed: 2026-01-03.

588 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-rag: Learning to retrieve, generate, and critique through self-reflection.

591 Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade, Vageesh D C, Arun Iyer, Suresh Parthasarathy, Sriram Rajamani, B. Ashok, and Shashank Shet. 2023. [Codeplan: Repository-level coding using LLMs and planning](#). In *NeurIPS 2023 Foundation Models for Decision Making Workshop*.

597 Amanda Bertsch, Maor Ivgi, Emily Xiao, Uri Alon, Jonathan Berant, Matthew R Gormley, and Graham Neubig. 2025. In-context learning with long-context models: An in-depth exploration. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 12119–12149.

605 Wei Cheng, Yuhan Wu, and Wei Hu. 2024. Dataflow-guided retrieval augmentation for repository-level code completion. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7957–7977.

610 Yangruibo Ding, Jinjun Peng, Marcus J. Min, Gail Kaiser, Junfeng Yang, and Baishakhi Ray. 2024a.

[Semcoder: Training code language models with comprehensive semantics reasoning](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 60275–60308. Curran Associates, Inc.

Yangruibo Ding, Zijian Wang, Wasi Ahmad, Hantian Ding, Ming Tan, Nihal Jain, Murali Krishna Ramanathan, Ramesh Nallapati, Parminder Bhatia, Dan Roth, and 1 others. 2023. Crosscodeeval: A diverse and multilingual benchmark for cross-file code completion. *Advances in Neural Information Processing Systems*, 36:46701–46723.

Yangruibo Ding, Zijian Wang, Wasi Ahmad, Murali Krishna Ramanathan, Ramesh Nallapati, Parminder Bhatia, Dan Roth, and Bing Xiang. 2024b. Cocomic: Code completion by jointly modeling in-file and cross-file context. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3433–3445.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2(1).

CRAC Generation. 2024. Coderag-bench: Can retrieval augment code generation.

Amirata Ghorbani and James Zou. 2019. Data shapley: Equitable valuation of data for machine learning. In *International conference on machine learning*, pages 2242–2251. PMLR.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, and 1 others. 2025. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.

Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, and 1 others. 2024. Qwen2.5-coder technical report. *CoRR*.

Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering. In *Proceedings of the 16th conference of the european chapter of the association for computational linguistics: main volume*, pages 874–880.

Paul Jaccard. 1912. The distribution of the flora in the alpine zone. 1. *New phytologist*, 11(2):37–50.

Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. In *Proceedings of the 2023*

667					
668					
669	Carlos E Jimenez, John Yang, Alexander Wettig,				
670	Shunyu Yao, Kexin Pei, Ofir Press, and Karthik				
671	Narasimhan. 2024. Swe-bench: Can language mod-				
672	els resolve real-world github issues? In <i>12th Inter-</i>				
673	<i>national Conference on Learning Representations,</i>				
674	<i>ICLR 2024.</i>				
675	Mintong Kang, Nezihe Merve Gürel, Ning Yu, Dawn				
676	Song, and Bo Li. 2024. C-RAG: Certified genera-				
677	tion risks for retrieval-augmented language models.				
678	In <i>Forty-first International Conference on Machine</i>				
679	<i>Learning.</i>				
680	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke				
681	Zettlemoyer, and Mike Lewis. 2020. Generalization				
682	through memorization: Nearest neighbor language				
683	models. In <i>International Conference on Learning</i>				
684	<i>Representations.</i>				
685	Denis Kocetkov, Raymond Li, Loubna Ben allal, Jia LI,				
686	Chenghao Mou, Yacine Jernite, Margaret Mitchell,				
687	Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf,				
688	Dzmitry Bahdanau, Leandro Von Werra, and Harm				
689	de Vries. 2023. The stack: 3 TB of permissively li-				
690	icensed source code. <i>Transactions on Machine Learn-</i>				
691	<i>ing Research.</i>				
692	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio				
693	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-				
694	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-				
695	täschel, and 1 others. 2020. Retrieval-augmented				
696	generation for knowledge-intensive nlp tasks. <i>Advances</i>				
697	<i>in neural information processing systems</i> , 33:9459–				
698	9474.				
699	R Li, LB Allal, Y Zi, N Muennighoff, D Kocetkov,				
700	C Mou, M Marone, C Akiki, J Li, J Chim, and 1				
701	others. 2023. Starcoder: May the source be with				
702	you! <i>Transactions on machine learning research.</i>				
703	Yanzhou Li, Shangqing Liu, Kangjie Chen, Tianwei				
704	Zhang, and Yang Liu. 2025. Impact-driven context				
705	filtering for cross-file code completion. In <i>Second</i>				
706	<i>Conference on Language Modeling.</i>				
707	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paran-				
708	jape, Michele Bevilacqua, Fabio Petroni, and Percy				
709	Liang. 2024a. Lost in the middle: How language				
710	models use long contexts. <i>Transactions of the Asso-</i>				
711	<i>ciation for Computational Linguistics</i> , 12:157–173.				
712	Tianyang Liu, Canwen Xu, and Julian McAuley. 2024b.				
713	Repobench: Benchmarking repository-level code				
714	auto-completion systems. In <i>The Twelfth Interna-</i>				
715	<i>tional Conference on Learning Representations.</i>				
716	Wei Liu, Ailun Yu, Daoguang Zan, Bo Shen, Wei Zhang,				
717	Haiyan Zhao, Zhi Jin, and Qianxiang Wang. 2024c.				
718	Graphcoder: Enhancing repository-level code				
719	completion via coarse-to-fine retrieval based on code				
720	context graph. In <i>Proceedings of the 39th IEEE/ACM</i>				
721	<i>International Conference on Automated Software En-</i>				
722	<i>gineering</i> , pages 570–581.				
	Scott M. Lundberg and Su-In Lee. 2017. A unified				723
	approach to interpreting model predictions. In <i>Pro-</i>				724
	<i>ceedings of the 31st International Conference on Neu-</i>				725
	<i>ral Information Processing Systems, NIPS’17</i> , page				726
	4768–4777, Red Hook, NY, USA. Curran Associates				727
	Inc.				728
	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,				729
	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.				730
	When not to trust language models: Investigating				731
	effectiveness of parametric and non-parametric				732
	memories. In <i>Proceedings of the 61st Annual Meeting of</i>				733
	<i>the Association for Computational Linguistics (Vol-</i>				734
	<i>ume 1: Long Papers)</i> , pages 9802–9822.				735
	Ikhtiyor Nematov, Tarik Kalai, Elizaveta Kuzmenko,				736
	Gabriele Fugagnoli, Dimitris Sacharidis, Katja Hose,				737
	and Tomer Sagi. 2025. Source attribution in				738
	retrieval-augmented generation. <i>arXiv preprint</i>				739
	<i>arXiv:2507.04480.</i>				740
	OpenAI. 2025. Gpt-5. https://openai.com/gpt-5 .				741
	Accessed: 2026-01-03.				742
	Md Rizwan Parvez, Wasi Ahmad, Saikat Chakraborty,				743
	Baishakhi Ray, and Kai-Wei Chang. 2021. Retrieval				744
	augmented code generation and summarization. In				745
	<i>Findings of the Association for Computational Lin-</i>				746
	<i>guistics: EMNLP 2021</i> , pages 2719–2734.				747
	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten				748
	Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi,				749
	Jingyu Liu, Romain Sauvestre, Tal Remez, and 1				750
	others. 2023. Code llama: Open foundation models				751
	for code. <i>arXiv preprint arXiv:2308.12950.</i>				752
	Lloyd S Shapley and 1 others. 1953. A value for n-				753
	person games.				754
	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan				755
	Scales, David Dohan, Ed H Chi, Nathanael Schärli,				756
	and Denny Zhou. 2023. Large language models can				757
	be easily distracted by irrelevant context. In <i>Inter-</i>				758
	<i>national Conference on Machine Learning</i> , pages				759
	31210–31227. PMLR.				760
	Disha Shrivastava, Denis Kocetkov, Harm De Vries,				761
	Dzmitry Bahdanau, and Torsten Scholak. 2023. Re-				762
	popofusion: Training code models to understand your				763
	repository. <i>arXiv preprint arXiv:2306.10998.</i>				764
	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017.				765
	Axiomatic attribution for deep networks. In <i>Interna-</i>				766
	<i>tional conference on machine learning</i> , pages 3319–				767
	3328. PMLR.				768
	Zhepei Wei, Wei-Lin Chen, and Yu Meng. 2025. In-				769
	structRAG: Instructing retrieval-augmented genera-				770
	tion via self-synthesized rationales. In <i>The Thirteenth</i>				771
	<i>International Conference on Learning Representa-</i>				772
	<i>tions.</i>				773
	Di Wu, Wasi Uddin Ahmad, Dejiao Zhang, Murali Kr-				774
	ishna Ramanathan, and Xiaofei Ma. 2024. Repo-				775
	former: selective retrieval for repository-level code				776

777 completion. In *Proceedings of the 41st International*
778 *Conference on Machine Learning, ICML'24*.
779 JMLR.org.

780 Yingtai Xiao, Yuqing Zhu, Sirat Samyoun, Wanrong
781 Zhang, Jiachen T Wang, and Jian Du. 2025. Token-
782 shapley: Token level context attribution with shapley
783 value. In *Findings of the Association for Computa-*
784 *tional Linguistics: ACL 2025*, pages 3882–3894.

785 Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2024. **RE-**
786 **COMP: Improving retrieval-augmented LMs with**
787 **context compression and selective augmentation**. In
788 *The Twelfth International Conference on Learning*
789 *Representations*.

790 Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling.
791 2024. Corrective retrieval augmented generation.
792 *arXiv preprint arXiv:2401.15884*.

793 Zezhou Yang, Ting Peng, Cuiyun Gao, Chaozheng
794 Wang, Hailiang Huang, and Yuetang Deng. 2025.
795 A deep dive into retrieval-augmented generation for
796 code completion: Experience on wechat. *arXiv*
797 *preprint arXiv:2507.18515*.

798 Chengyuan Yao and Satoshi Fujita. 2024. Adaptive
799 control of retrieval-augmented generation for large
800 language models through reflective tags. *Electronics*,
801 13(23):4643.

802 Zikun Ye and Hema Yoganasimhan. 2025. Fair docu-
803 ment valuation in llm summaries via shapley values.
804 *arXiv preprint arXiv:2505.23842*.

805 Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Be-
806 rant. Making retrieval-augmented language models
807 robust to irrelevant context. In *ICLR 2024 Workshop*
808 *on Large Language Model (LLM) Agents*.

809 Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin
810 Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao
811 Zeng, Jiajie Zhang, and 1 others. 2025. G1m-4.5:
812 Agentic, reasoning, and coding (arc) foundation mod-
813 els. *arXiv preprint arXiv:2508.06471*.

814 Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin
815 Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and
816 Weizhu Chen. 2023. Repocoder: Repository-level
817 code completion through iterative retrieval and gener-
818 ation. In *The 2023 Conference on Empirical Methods*
819 *in Natural Language Processing*.

820 Sheng Zhang, Yifan Ding, Shuquan Lian, Shun Song,
821 and Hui Li. 2025. Coderag: Finding relevant and nec-
822 cessary knowledge for retrieval-augmented repository-
823 level code completion. In *Proceedings of the 2025*
824 *Conference on Empirical Methods in Natural Lan-*
825 *guage Processing*, pages 23289–23299.

A Details of dataset construction

A.1 Metrics Formulation

We evaluate code completion quality using Exact Match (EM), Edit Similarity (ES), and Unit Tests (UT). Let \hat{Y} be the generated code and Y be the ground truth:

$$\begin{aligned} \text{EM} &= \mathbb{1}(\hat{Y} = Y) \\ \text{ES} &= \left[1 - \frac{\mathcal{D}(\hat{Y}, Y)}{\max(|\hat{Y}|, |Y|)}\right] \times 100\% \\ \text{UT} &= \mathbb{1}(\text{Pass}(\hat{Y})) \end{aligned}$$

where $\mathbb{1}(\cdot)$ is the indicator function, $\mathcal{D}(\cdot)$ denotes the Levenshtein distance, and $\text{Pass}(\cdot)$ returns true if the code passes all unit tests.

A.2 Data Collection and Preprocessing

File-level filtering. We begin with conservative file hygiene to reduce retrieval noise and stabilize likelihood-based labeling. We keep only .py files and discard files with fewer than 10 non-empty lines. To remove minified/generated blobs that distort sparse retrieval, we drop files whose maximum line length exceeds 300 characters or whose average line length exceeds 120 characters (computed after trimming trailing whitespace). We further filter out non-code payloads by requiring alphanumeric density ≥ 0.35 (ratio of letters/digits over all characters). Finally, we exclude vendored or generated directories by path keywords, including vendor/, third_party/, site-packages/, dist/, build/, .venv/, and migrations/. All statistics are computed on UTF-8 decoded text (with a permissive fallback that drops undecodable bytes).

Repository-level filtering. We retain repositories with sufficient structure for cross-file interactions by requiring at least 8 remaining Python files and total non-empty LOC between 300 and 50,000 after file-level filtering. To avoid duplicate-heavy projects where top- K retrieval collapses to repeated copies, we estimate the near-duplicate file ratio using **SimHash**. Specifically, for each repository we compute SimHash fingerprints over a normalized UTF-8 text representation of each file (whitespace-collapsed, with trailing whitespace removed), and perform pairwise checks on up to the first 200 files (`max_files_for_dup_check=200`). We mark two files as near-duplicates if

their SimHash Hamming distance is at most 3 (`simhash_hamming_threshold=3`). Repositories with more than 30% near-duplicate files are discarded (`max_dup_ratio=0.3`).

We also enforce syntactic integrity by parsing a sampled subset of files with `Python ast.parse`. Concretely, we uniformly sample up to 20 files per repository (`ast_sample_k=20`) from the remaining Python files after file-level filtering, and compute the parse success rate on this sample. Repositories with parse success rate $< 70\%$ are removed (`min_ast_parse_rate=0.7`). For reproducibility, both the duplicate-check subsampling (when applicable) and AST sampling use a fixed random seed of 13 (`seed=13`).

A.3 Data labeling.

We chunk the cross-file pool and construct a retrieval query from the in-file context for each target span Y . For each query, we retrieve the top- K candidate chunks and distill coalition-aware decisions into chunk-wise KEEP/DROP labels and a retrieval-control token. Specifically, given $X_{\text{in}} = (X_p, X_s)$ and retrieved candidates $X_{\text{cc}} = (cc_1, \dots, cc_K)$, we run our offline **ChunkShapley** pipeline to obtain a verification-selected coalition $S^* \subseteq \{1, \dots, K\}$. We first compute a teacher-forced baseline log-likelihood $\ell(\emptyset)$ and probe each chunk in isolation, $\Delta_i = \ell(\{i\}) - \ell(\emptyset)$, yielding a signed vote $y_i = \text{sign}(\Delta_i)$ and weight $\omega_i = |\Delta_i|$. We then define a lightweight surrogate game $v_{\text{sur}}(S) = \sigma(\beta \sum_{i \in S} \omega_i y_i) - \sigma(0)$, where $\sigma(\cdot)$ is the sigmoid function, and compute exact surrogate Shapley values by enumerating all 2^K coalitions (tractable for small K and performed offline). Finally, we verify a bounded set of Shapley-proposed coalitions using the frozen generator under decoding-time constraints and select S^* that maximizes completion quality (lexicographically by ES then EM). We treat S^* as the teacher subset and assign labels: $Q(cc_i) = \langle \text{KEEP} \rangle$ if $i \in S^*$ and $\langle \text{DROP} \rangle$ otherwise.

To supervise retrieval triggering, we assign the retrieval-control token $r^* \in \{\langle \text{NEED} \rangle, \langle \text{DONE} \rangle\}$ by comparing the completion quality with and without cross-file evidence: if the in-file context alone already achieves ES above a threshold τ_{done} (or the verified coalition provides negligible gain), we set $r^* = \langle \text{DONE} \rangle$; otherwise $r^* = \langle \text{NEED} \rangle$. To ensure that retained evidence is meaningful, we filter instances by requiring the verification-selected coalition to achieve $\text{ES} \geq \tau_{\text{es}}$. Alg. 2 summarizes the labeling procedure.

Table 3: Code completion performance in the Left-to-Right setting.

Model	Strategy	RepoEval						CCLongEval			CCEval	
		Line		API		Function		Chunk		Func	Line	
		EM	ES	EM	ES	UT	ES	EM	ES	ES	EM	ES
SC-Base-1B	No-Retrieve	33.42	57.88	28.54	57.36	16.55	40.21	22.45	53.05	39.88	16.54	54.47
	Full-Retrieve	44.52	66.21	36.95	64.77	21.30	48.55	31.12	63.49	45.36	20.12	58.21
	RepoFormer	46.12	68.33	37.44	66.12	23.45	50.12	32.55	65.12	46.88	22.45	60.33
	CODEFILTER	48.88	70.15	39.85	69.11	24.12	51.55	34.15	66.88	48.22	24.88	62.55
	REPOSHAPLEY	54.21 ^{+5.33}	76.45 ^{+6.30}	45.66 ^{+5.81}	74.88 ^{+5.77}	29.85 ^{+5.73}	57.22 ^{+5.67}	40.55 ^{+6.40}	72.44 ^{+5.56}	54.12 ^{+5.90}	30.12 ^{+5.24}	68.95 ^{+6.40}
SC-Base-3B	No-Retrieve	35.82	60.12	29.55	58.45	20.05	38.95	25.44	53.45	44.82	18.22	57.51
	Full-Retrieve	50.45	70.88	40.66	67.89	26.12	48.66	36.15	59.88	45.75	23.45	62.12
	RepoFormer	50.11	71.95	41.02	69.88	27.55	50.12	36.75	61.95	47.12	25.66	63.88
	CODEFILTER	53.12	73.66	43.15	73.12	27.88	51.22	38.05	61.55	48.88	27.45	65.12
	REPOSHAPLEY	59.45 ^{+6.33}	79.11 ^{+5.45}	48.88 ^{+5.73}	79.55 ^{+6.43}	33.45 ^{+5.57}	57.88 ^{+6.66}	44.22 ^{+6.17}	67.12 ^{+5.57}	54.66 ^{+5.78}	33.15 ^{+5.70}	70.44 ^{+5.32}
SC-Base-7B	No-Retrieve	38.15	62.12	31.45	59.88	21.88	39.95	29.45	58.55	53.45	19.68	59.00
	Full-Retrieve	51.22	71.55	42.45	68.33	28.15	50.45	41.55	64.88	48.75	24.55	63.45
	RepoFormer	50.88	70.45	40.88	72.66	28.05	48.55	41.05	64.75	48.15	26.88	65.12
	CODEFILTER	53.95	74.22	44.55	72.15	29.05	52.33	42.12	66.88	57.88	28.55	67.55
	REPOSHAPLEY	59.88 ^{+5.93}	80.55 ^{+6.33}	50.12 ^{+5.57}	78.45 ^{+6.30}	34.66 ^{+5.61}	58.12 ^{+5.79}	48.45 ^{+6.33}	72.15 ^{+5.27}	63.45 ^{+5.57}	34.12 ^{+5.72}	73.22 ^{+5.67}
Llama-7B	No-Retrieve	39.55	64.12	30.88	60.22	22.45	42.55	30.12	58.12	45.45	20.88	60.12
	Full-Retrieve	52.45	70.88	43.15	68.75	26.88	50.12	41.22	63.88	52.66	25.44	64.55
	RepoFormer	51.12	71.45	40.88	70.88	28.66	50.05	39.45	62.95	50.45	27.12	65.88
	CODEFILTER	53.66	73.12	43.88	73.55	29.88	50.88	41.88	65.45	53.55	29.45	67.88
	REPOSHAPLEY	59.12 ^{+5.46}	78.66 ^{+5.54}	49.55 ^{+5.67}	79.12 ^{+5.57}	35.12 ^{+5.24}	56.45 ^{+5.57}	47.22 ^{+5.34}	71.05 ^{+5.60}	59.12 ^{+5.57}	35.66 ^{+6.21}	73.45 ^{+5.57}
Llama-13B	No-Retrieve	41.55	65.12	31.22	60.66	24.12	43.55	31.55	57.66	46.12	21.88	61.45
	Full-Retrieve	54.22	74.45	44.66	71.88	28.95	51.45	43.22	68.12	50.22	26.55	66.12
	RepoFormer	52.45	71.55	43.95	71.66	28.66	51.12	43.55	67.88	52.45	28.12	67.45
	CODEFILTER	55.33	75.12	45.12	74.95	30.12	52.45	44.22	67.95	57.12	30.88	69.55
	REPOSHAPLEY	61.12 ^{+5.79}	81.45 ^{+6.33}	51.55 ^{+6.43}	80.66 ^{+5.71}	36.45 ^{+6.33}	58.22 ^{+5.77}	50.12 ^{+5.90}	73.45 ^{+5.50}	62.88 ^{+5.76}	36.95 ^{+6.07}	75.12 ^{+5.57}

A.4 Labeling Algorithm Details

Algorithm 3 outlines the process of deriving supervision signals from raw code repositories. First, top- K candidate chunks X_{cc} are retrieved based on the query window Q . We then utilize CHUNKSHAPLEY to identify the optimal chunk subset S^* that maximizes generation quality relative to the ground truth Y .

The labeling logic follows three specific criteria:

- Quality Control:** Instances are discarded if the optimal subset’s performance falls below a minimum threshold τ_{es} , ensuring training data quality.
- Retrieval Label (r^*):** We measure the performance gain of using external context (S^*) versus the closed-book baseline (\emptyset). If the gain is negligible ($\leq \epsilon$), the retrieval label is set to $\langle \text{DONE} \rangle$; otherwise, it is $\langle \text{NEED} \rangle$.
- Selection Label (q_i^*):** Individual chunks are labeled as $\langle \text{KEEP} \rangle$ if they belong to the optimal subset S^* , and $\langle \text{DROP} \rangle$ otherwise.

B Hyperparameter Optimization

We tune training hyperparameters using **StarCoderBase-1B** as a proxy model to reduce search cost. Unless otherwise specified, all other settings follow the main experimental setup (e.g., data split, prompt formats, max sequence length, and batching).

Search space. We conduct a grid search on the following space: learning rate $\in \{1 \times 10^{-5}, 2 \times 10^{-5}, 5 \times 10^{-5}\}$, loss weight $\lambda \in \{0.2, 1.0, 2.0, 5.0\}$, training epochs $\in \{1, 2, 5\}$, and warmup steps $\in \{50, 100\}$. Here λ is applied to the retrieval-control and selection losses, i.e., $\lambda_R = \lambda_S = \lambda$ (while the generation loss uses unit weight).

Selection criterion. For each configuration, we evaluate code completion performance on the validation split using the same metrics as in the main experiments. We select the best hyperparameters by maximizing the validation completion quality (with ES as the primary criterion and EM as a tie-breaker).

Final configuration. The selected hyperparameters are: learning rate 2×10^{-5} , $\lambda_R = \lambda_S = 2.0$, epochs = 2, warmup steps = 50. We reuse this configuration for all backbones in our experiments for consistency.

C Detailed Experiments

C.1 Ablation on Verification Scope (L)

Table 5 shows that expanding the verification scope from $L = 0$ to $L = 3$ brings the largest gains: moving beyond prefix-only verification substantially improves both EM and ES, and performance increases steadily up to the default $L = 3$. In contrast, further enlarging the scope ($L > 3$) yields only marginal improvements, despite a rapidly growing

Algorithm 2: ChunkShapley: Surrogate Shapley Attribution with Bounded Verification

Input: In-file context X_{in} ; ground-truth completion Y ; retrieved chunks $X_{\text{cc}} = (cc_1, \dots, cc_K)$; Frozen generator G_θ ; surrogate scale β ; verification params (N_v, L) .
Output: Verification-selected coalition $S^* \subseteq \{1, \dots, K\}$; surrogate Shapley scores $\{\phi_i\}_{i=1}^K$.

```

1  $\ell(\emptyset) \leftarrow \frac{1}{|Y|} \log p_\theta(Y | X_{\text{in}})$ 
2 for  $i \leftarrow 1$  to  $K$  do
3    $\ell(\{i\}) \leftarrow \frac{1}{|Y|} \log p_\theta(Y | X_{\text{in}}, \{cc_i\})$ 
4    $\Delta_i \leftarrow \ell(\{i\}) - \ell(\emptyset)$ 
5    $y_i \leftarrow \text{sign}(\Delta_i)$ ;  $\omega_i \leftarrow |\Delta_i|$ 
6 end
7 foreach  $S \subseteq \{1, \dots, K\}$  do
8    $g(S) \leftarrow \sum_{j \in S} \omega_j y_j$ 
9    $v_{\text{sur}}(S) \leftarrow \sigma(\beta g(S)) - \sigma(0)$ 
10 end
11 for  $i \leftarrow 1$  to  $K$  do
12    $\phi_i \leftarrow 0$ 
13   foreach  $S \subseteq \{1, \dots, K\} \setminus \{i\}$  do
14      $w(S) \leftarrow \frac{|S|!(K-|S|-1)!}{K!}$ 
15      $\phi_i \leftarrow \phi_i + w(S)(v_{\text{sur}}(S \cup \{i\}) - v_{\text{sur}}(S))$ 
16   end
17 end
18  $\pi_\phi \leftarrow \text{argsort}(\{\phi_i\}, \text{desc})$ ;
19  $\pi_\Delta \leftarrow \text{argsort}(\{\Delta_i\}, \text{desc})$ 
20  $\mathcal{C} \leftarrow \text{BuildPool}(\pi_\phi, \pi_\Delta; N_v, L)$ 
21 foreach  $S \in \mathcal{C}$  do
22    $\hat{Y}_S \leftarrow \text{Decode}(G_\theta | X_{\text{in}}, X_S)$ 
23   Compute ES( $\hat{Y}_S, Y$ ) and EM( $\hat{Y}_S, Y$ )
24 end
25  $S^* \leftarrow \arg \max_{S \in \mathcal{C}} (\text{ES}(\hat{Y}_S, Y), \text{EM}(\hat{Y}_S, Y))$ 
26 return  $S^*$ ,  $\{\phi_i\}_{i=1}^K$ 

```

979 candidate pool and offline labeling cost. Therefore,
980 we adopt $L = 3$ as a practical default that captures
981 most of the benefit of combinatorial probing.

982 C.2 Oracle Analysis

983 To validate the theoretical superiority of Shapley-
984 based valuation over independent likelihood prob-
985 ing (as used in CODEFILTER), we conducted an Or-
986 acle study. We calculated the best possible Edit
987 Similarity (ES) achievable if the model perfectly
988 selected chunks according to the respective valua-
989 tion methods (selecting top- K chunks with score
990 > 0).

991 As shown in Table 6, the Shapley-based oracle
992 outperforms the CODEFILTER oracle by **10.45** per-
993 centage points. This confirms that modelling chunk
994 interactions, like synergy and conflict, is critical for
995 repository-level code completion, as independent
996 probing fails to identify chunks that are only useful

Algorithm 3: REPOSHAPLEY Cross-file Labeling via ChunkShapley

Input: Repository cross-file pool X_{out} ; in-file context $X_{\text{in}} = (X_p, X_s)$; target span Y ; Retriever R ; frozen generator G ; chunk window w ; stride s ; retrieve budget K ; verification params (N_v, L) (as in Alg. 2); thresholds τ_{es} and ϵ .
Output: Labeled instance: retrieval label $r^* \in \{\langle \text{NEED} \rangle, \langle \text{DONE} \rangle\}$ and selection labels (q_1^*, \dots, q_K^*) with $q_i^* \in \{\langle \text{KEEP} \rangle, \langle \text{DROP} \rangle\}$

```

1  $Q \leftarrow X_p[-w : ]$ 
2  $\tilde{X}_{\text{out}} \leftarrow \text{chunkize}(X_{\text{out}}; w, s)$ 
3  $X_{\text{cc}} \leftarrow R(Q, \tilde{X}_{\text{out}})[1:K]$ 
4  $(S^*, \hat{Y}_{S^*}) \leftarrow \text{ChunkShapley}(X_{\text{in}}, Y, X_{\text{cc}}, G; N_v, L)$ 
5  $\hat{Y}_\emptyset \leftarrow G(X_{\text{in}})$ 
6 if  $\text{ES}(\hat{Y}_{S^*}, Y) < \tau_{\text{es}}$  then
7   return discard instance
8 end
9 if  $\text{ES}(\hat{Y}_{S^*}, Y) - \text{ES}(\hat{Y}_\emptyset, Y) \leq \epsilon$  then
10    $r^* \leftarrow \langle \text{DONE} \rangle$ 
11 end
12  $r^* \leftarrow \langle \text{NEED} \rangle$ 
13 for  $i \leftarrow 1$  to  $K$  do
14    $q_i^* \leftarrow \langle \text{KEEP} \rangle$  if  $i \notin S^*$  then
15      $q_i^* \leftarrow \langle \text{DROP} \rangle$ 
16   end
17 end
18 end
19 return  $r^*$ ,  $(q_1^*, \dots, q_K^*)$ 

```

Table 4: Hyperparameter search space and selected values (tuned on StarCoderBase-1B).

Hyperparameter	Search space	Selected
Learning rate	{1e-5, 2e-5, 5e-5}	2e-5
λ ($\lambda_R = \lambda_S$)	{0.2, 1.0, 2.0, 5.0}	2.0
Epochs	{1, 2, 5}	2
Warmup steps	{50, 100}	50

when combined, for instance interface definitions
and implementations.

C.3 Sensitivity Analysis on Retrieval Budget K

Table 7 investigates the trade-off between comple-
tion performance and inference latency by varying
the retrieval budget K (i.e., the number of candi-
date chunks processed by ChunkShapley) on SC-
Base-1B. Increasing K expands the search space
for complementary evidence, potentially captur-
ing more synergistic interactions. However, since
our method involves exact Shapley estimation via
subset enumeration, the computational cost grows
exponentially with K . Specifically, a small budget
($K = 7$) yields low latency but fails to retrieve suffi-

Table 5: **Ablation on Verification Scope (L)**. Impact of the Top- L range used for combinatorial probing in the post-verification stage. We report the average size of the candidate pool $|\mathcal{C}|$, offline labeling latency per instance, and performance on SC-Base-1B.

Top- L	Labeling Cost		Performance	
	Avg. Pool Size ($ \mathcal{C} $)	Train Time per Sample (s)	EM (%)	ES (%)
0 (Prefix Only)	8.41	33	45.15	70.86
1	10.78	94	57.73	76.27
2	12.56	187	59.15	78.40
3(Default)	18.29	348	61.34	82.78
4	25.07	671	61.70	83.22
5	35.42	869	61.94	83.24
7	68.94	1528	62.13	83.59
10 (All)	225	6823	-	-

Table 6: **Oracle Performance Comparison**. We report the mean Best ES score achievable by selecting contexts based on oracle labels. REPOSHAPLEY (Oracle) demonstrates a significantly higher theoretical upper bound.

Method	Oracle Best ES (%)
Full-Retrieve	71.52
CODEFILTER (Oracle)	85.23
REPOSHAPLEY (ORACLE)	95.68

cient complementary pairs, resulting in suboptimal accuracy (52.16% EM). Conversely, increasing K beyond 10 yields *diminishing returns*; for instance, expanding to $K = 13$ marginally improves ES by 0.51% but causes latency to explode to over 3.5 seconds due to the combinatorial complexity of the surrogate game (2^{13} subsets), rendering it impractical. Crucially, $K = 10$ achieves the optimal balance, we adopt $K = 10$ as the default setting to balance interaction coverage with inference efficiency.

C.4 Ablation Study on Coalition Utility Functions

In **ChunkShapley**, the choice of the characteristic function $v(S)$ is critical, as it defines the "value" distributed among retrieved chunks. We hypothesize that while task-specific metrics (like Exact Match) align perfectly with the final objective, they provide sparse and noisy signals for attribution. To validate the effectiveness of our Log-likelihood-based utility, we conduct an ablation study comparing it against task-metric-based utilities.

Table 7: Sensitivity analysis of the retrieval budget (K) on SC-Base-1B in the Infilling setting. We report Line Completion accuracy (EM/ES) and average inference latency. $K = 10$ achieves the best trade-off between context coverage and computational cost.

Retrieval Size	RepoEval-Line		Efficiency
	EM	ES	Latency (ms)
7	52.16	70.44	513
9	58.33	75.12	833
10 (Ours)	61.34	82.78	1053
11	61.37	82.40	1924
13	61.88	81.62	3539
20	61.76	79.92	18825

Experimental Setup. We compare three definitions of coalition utility $v(S)$:

Log-likelihood Utility (Ours): We use the normalized token-level log-probability gain under teacher forcing. This provides a continuous, dense signal reflecting the model’s confidence:

$$v_{\log}(S) = \frac{1}{|Y|} \sum_{t=1}^{|Y|} \left(\log p_{\theta}(y_t | y_{<t}, X_{\text{in}}, X_S) - \log p_{\theta}(y_t | y_{<t}, X_{\text{in}}, \emptyset) \right).$$

Exact Match (EM) Utility: We define utility as the binary gain in obtaining a perfect prediction. This signal is discrete ($\in \{-1, 0, 1\}$) and highly sparse:

$$v_{\text{EM}}(S) = \mathbb{1}[\text{EM}(\hat{Y}_S, Y) = 1] - \mathbb{1}[\text{EM}(\hat{Y}_{\emptyset}, Y) = 1]$$

Table 8: **Ablation study on Utility Functions.** We report the performance of **REPOSHAPLEY** when trained with Shapley labels derived from different utility definitions. *Log-likelihood* (Ours) significantly outperforms metric-based utilities due to the density and stability of the teacher-forcing signal.

Utility Function	Signal Type	EM	ES
w/ v_{EM}	Binary	58.12	77.45
w/ v_{ES}	Discrete-Step	59.03	78.10
w/ v_{log} (Ours)	Continuous	60.50	79.07

where \hat{Y}_S is the greedy decoding result given context S .

Edit Similarity (ES) Utility: We define utility based on the improvement in surface-level similarity. While continuous (0-100), ES is derived from discrete decoding steps and is non-differentiable:

$$v_{ES}(S) = ES(\hat{Y}_S, Y) - ES(\hat{Y}_\emptyset, Y) \quad (1)$$

For all variants, we compute the exact Shapley values using the respective $v(S)$, select the optimal subset S^* using the verification strategy, and train the corresponding **REPOSHAPLEY** model.

Results and Analysis. As shown in Table 8, using log-likelihood as the utility function yields the best performance. We attribute the inferiority of metric-based utilities (v_{EM}, v_{ES}) to the **sparsity and high variance** of the signal. In code completion, a chunk might significantly improve the model’s understanding (providing the correct variable type) without immediately flipping the final prediction to an exact match. v_{log} captures this "partial credit," whereas v_{EM} assigns zero value, leading to false negatives in attribution. Furthermore, generation-based metrics are sensitive to decoding dynamics, where a small change in context might drastically alter the greedy path, causing $v_{metric}(S)$ to fluctuate wildly. In contrast, teacher-forced log-probabilities provide a smoother and more robust estimation of marginal contribution.

C.5 Filtering effectiveness via selective drop and counterfactual inverse.

Following **CODEFILTER** (Li et al., 2025), we study whether attribution signals can reliably separate helpful from harmful retrieved context on RepoEval under the same Jaccard-based retriever. Table 9 reports two complementary interventions: **Selective (Drop)**, which removes chunks labeled

Table 9: Impact of the filtering policy based on different attribution signals. **Selective** denotes removing chunks labeled as $\langle \text{DROP} \rangle$, while **Inverse** retains only those chunks (to verify the toxicity of dropped content).

Strategy	CODEFILTER		REPOSHAPLEY	
	EM	ES	EM	ES
Selective (Drop)	49.31	57.50	54.79\uparrow	78.62\uparrow
Inverse (Keep)	33.17	40.26	34.72	41.15

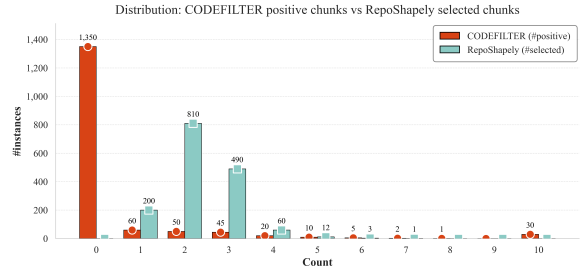


Figure 4: Distribution: CODEFILTER positive chunks vs. REPOSHAPLEY selected chunks

as $\langle \text{DROP} \rangle$ and keeps the remaining context, and **Inverse (Keep)**, which retains only the dropped chunks as a counterfactual diagnostic.

REPOSHAPLEY gains markedly from selective filtering, improving both EM and ES compared to the **CODEFILTER** counterpart. In contrast, the inverse setting substantially degrades performance for both methods, confirming that the dropped chunks are predominantly low-utility (e.g., redundant or misleading) rather than accidentally filtered-out evidence. Notably, Figure 4 shows that **CODEFILTER**’s decisions are prone to brittle single-chunk thresholds: when individual signals are weak, it tends to label very few chunks as positive, effectively collapsing the available context. **REPOSHAPLEY** instead maintains a stable selection set, consistent with interaction-aware supervision that removes toxic context while preserving the evidence required for accurate repository-level completion.

C.6 Interaction-Aware Proposal.

We examine whether gains stem solely from verification by testing **Delta+Verify**, which replaces Shapley candidates with single-chunk rankings (Δ_i) under the same verifier. Table 10 shows **REPOSHAPLEY** achieves a higher net gain. This confirms that Δ scores miss **synergistic chunks** (weak in isolation), creating a hard performance ceiling. In contrast, Shapley effectively captures these high-

Table 10: **Impact of Proposal Mechanism.** Comparison between using Single-Chunk Δ vs. Coalition Shapley ϕ to generate candidates for the verification step.

Proposal Method	Metric	Gain over Full
Delta (Δ) + Verify	ES	74.21
REPOSHAPLEY (ϕ) + Verify	ES	82.78

potential interactive subsets, providing the verifier with a superior candidate pool.

C.7 Abstention and Selection Analysis

To further understand RepoShapley’s decision behavior, we conduct an abstention analysis across datasets and task granularities. We partition retrieved candidates into three outcome types by combining the retriever score (high vs. low) with RepoShapley’s keep or drop decision: **High retained** (high-score kept), **High discarded** (high-score dropped), and **Low captured** (low-score kept). Figure 5 reports the proportions of these categories.

Across relatively local settings, RepoShapley stays consistent with the dense retriever while making selective corrections. On **RepoEval Line** and **CCLEval Chunk**, most high-scoring chunks are retained (93% and 90%), yet the policy still rejects a small fraction of high-score candidates (3–4%) and recovers a small but non-trivial set of low-score evidence (4–6%). This indicates that RepoShapley does not simply follow similarity ranking; instead, it can stably abstain from a few high-score but ineffective chunks and preserve occasional low-score signals even when surface similarity is largely reliable.

As task difficulty increases, the policy increasingly captures low-score but necessary evidence. For **RepoEval Function** and **CCLEval Function**, the low-captured portion rises to 12% and 11%, while the high-discarded portion remains modest (6% and 5%). The shift is most pronounced on **RepoEval API**, where RepoShapley retains only 69% of high-score chunks, discards 7% as redundant/noisy, and captures 24% from the low-score tail. This trend supports our core claim: when cross-file interactions become more complex, RepoShapley can recover individually under-ranked chunks whose utility emerges through strong interactions with other context, effectively filtering “high-score noise” while identifying “low-score but interaction-critical” evidence.

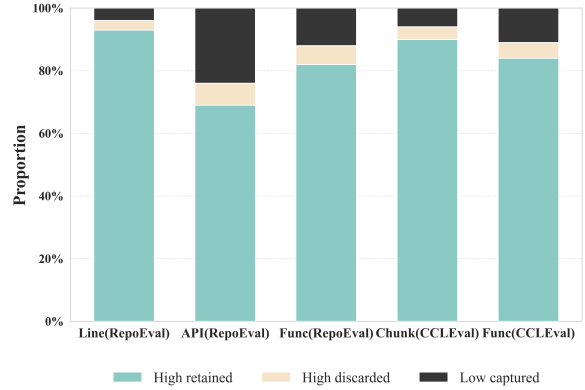


Figure 5: Breakdown of chunk selection decisions by RepoShapley on different benchmarks. **High retained** indicates consensus between retriever and policy. **High discarded** and **Low captured** highlight cases where RepoShapley corrects the retriever’s judgment.

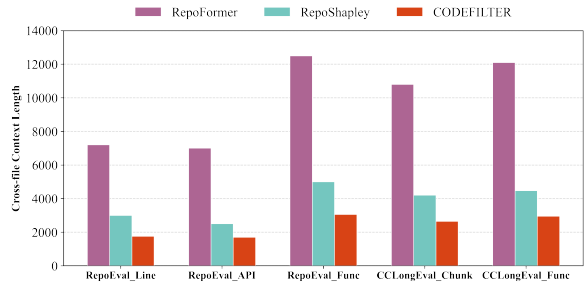


Figure 6: Comparison of retained cross-file context lengths. RepoFormer keeps the most tokens, CODEFILTER is the most aggressive, and RepoShapley balances pruning and coverage.

C.8 Context Length Distribution Analysis

Beyond accuracy, we analyze efficiency by measuring how many cross-file tokens are retained after filtering. Figure 6 shows the distribution of retained context lengths (in tokens) for RepoFormer, CODEFILTER, and RepoShapley across five benchmarks.

RepoFormer consistently retains the longest contexts and can exceed 12k tokens on function-level tasks, indicating a high-recall behavior that carries substantial redundancy. This increases both computational cost and the risk of distracting evidence during generation.

CODEFILTER sits at the other extreme: it retains the fewest tokens by pruning aggressively, which reduces latency but can remove weak yet necessary dependencies.

RepoShapley lies between these two regimes. It substantially shortens contexts relative to RepoFormer, suggesting effective removal of redun-

Table 11: Accuracy of state-of-the-art code LMs as the generation model and with REPOSHAPLEY as the policy model for selective RAG in Infilling setting. We compare DeepSeek (Guo et al., 2025), Qwen (Hui et al., 2024), GLM (Zeng et al., 2025), GPT-5 (OpenAI, 2025), and Claude Opus (Anthropic, 2025). For closed-source models, we use the official APIs for generation.

Model	Strategy	RepoEval-Line		RepoEval-API	
		EM	ES	EM	ES
StarCoderBase-7B	Full-Retrieve	58.26	77.79	50.38	75.01
	RepoShapley	65.81	86.59	58.79	84.11
CodeLlama-13B	Full-Retrieve	52.73	71.91	42.28	69.57
	RepoShapley	68.89	87.11	57.66	83.41
DeepSeek R1	Full-Retrieve	52.94	70.09	42.61	70.93
	RepoShapley	68.79	87.92	58.96	84.07
Qwen 2.5 Max	Full-Retrieve	53.38	71.91	44.72	72.39
	RepoShapley	69.51	88.17	59.14	84.58
GLM-4.5	Full-Retrieve	55.93	73.04	45.27	73.21
	RepoShapley	69.96	88.49	60.32	84.94
ChatGPT-5	Full-Retrieve	56.51	73.73	45.94	73.82
	RepoShapley	70.37	89.08	60.19	84.41
Claude Opus 4.1	Full-Retrieve	58.23	74.38	47.56	75.61
	RepoShapley	71.09	89.13	61.74	85.92

dant or conflicting chunks, while remaining slightly longer than CODEFILTER. This gap is expected: RepoShapley tends to keep complementary chunks that may appear low-signal in isolation but improve the coalition, consistent with the abstention analysis in Appendix C.7. Overall, RepoShapley achieves a favorable trade-off by reducing token overhead without sacrificing semantic coverage.

C.9 SOTA Generation Models with REPOSHAPLEY Policy in Infilling

Table 11 reports results when we pair REPOSHAPLEY with a wide range of state-of-the-art code LMs under the infilling setting on RepoEval-Line and RepoEval-API. For each backbone, we compare a standard *Full-Retrieve* strategy against using REPOSHAPLEY as a selective RAG policy while keeping the same generation model.

Across all backbones, REPOSHAPLEY consistently improves both EM and ES on both benchmarks, indicating that coalition-aware evidence filtering is complementary to model scaling and remains effective for both open-source and closed-source generators. For closed-source models, we use the official APIs for generation, while RE-

POSHAPLEY is applied as an external policy to decide whether to retrieve and which chunks to keep.

C.10 Latency-Accuracy Trade-off Analysis

Following RepoFormer (Wu et al., 2024), we visualize the latency-accuracy trade-off to evaluate the efficiency of REPOSHAPLEY across StarCoderBase-1B, 3B, and 7B as shown in Figure 7-9. By varying the retrieval triggering threshold t_c during inference, we control the model’s sensitivity to external evidence.

The results demonstrate that REPOSHAPLEY establishes a superior Pareto frontier compared to static retrieval strategies. We find that our model can improve accuracy while also reducing latency by skipping retrieval when the in-file context is already sufficient, and focusing retrieval on harder cases that truly need cross-file information. Consistent with prior observations, this efficiency gain is particularly pronounced in Line and API completion tasks, where avoiding the overhead of unnecessary retrieval significantly lowers average latency without compromising generation quality.

D Case Study

In this section, we present a case study to illustrate how REPOSHAPLEY performs interaction-aware chunk selection for repository-level code completion in the FIM setting. The target file defines utilities for extracting event start/end markers from log search results and computing event durations. In this instance, the missing span lies inside `LogEventStats.run`: after parsing each end-marker timestamp end from `end_tag` results, the code should immediately register it into `EventCollection` via `add_event_end`. After registering each end marker, the routine proceeds to handle start markers and finally calls `calculate_event_deltas`. The repository contains many timestamp-related helpers; however, most retrieved evidence is only partially relevant or unrelated (e.g., YAML formatting, tests, Ceph helpers). Naively appending all retrieved contexts can distract the model into rewriting timestamp parsing rather than emitting the required event-collection logic.

Instance (FIM). Given the in-file prefix and suffix (Figure 10), the model must generate the missing span at `<MID>`. Concretely, the correct completion should insert an `add_event_end` call that uses the current `result`'s `event_id` and the parsed end timestamp. For compact presentation, Figure 10 splits the code across two columns. Therefore, the `<SFX>` panel shows the subsequent lines after the insertion point (not necessarily the immediate next line in the source file), while preserving the original indentation and control flow.

Retrieved top-10 cross-file chunks. We retrieve the top-10 candidates $\{c_1, \dots, c_{10}\}$ from other files in the same repository (Figure 11). Most candidates are either unrelated or only partially relevant. Importantly, the most helpful timestamp utilities are split across multiple chunks (c_1, c_8, c_9), so that no single chunk alone fully specifies the needed behavior, making the evidence *interaction-heavy*.

Why the kept coalition matters. Chunks c_1, c_8 , and c_9 form a coherent timestamp-handling subroutine. c_1 and c_8 provide compatible `datetime` parsing formats, and c_9 implements temporal filtering logic used by the surrounding utilities. Full-Retrieve is distracted by irrelevant utilities (e.g., c_2, c_3) and hallucinates redundant timestamp-parsing logic inside `run`, instead of emitting the required `add_event_end` registration.

Generation comparison. Figure 12 compares the generations. REPOSHAPLEY correctly inserts the `add_event_end` registration and then follows the existing start-marker logic, whereas Full-Retrieve is distracted and produces redundant parsing code.

1271
1272
1273
1274
1275
1276

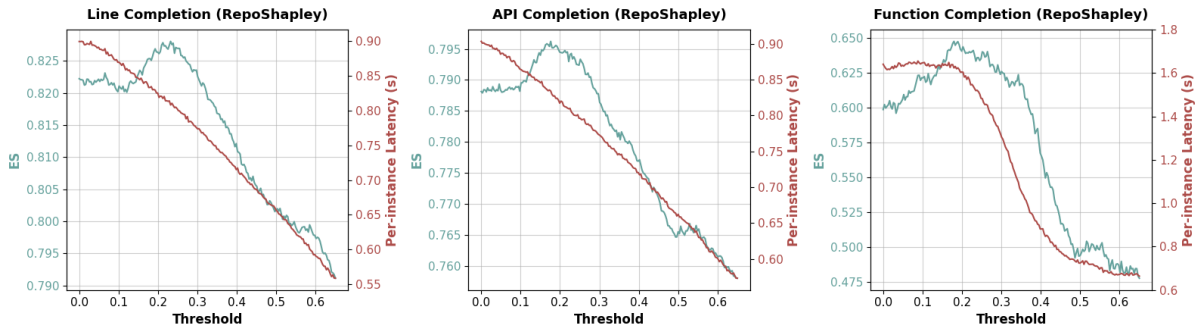


Figure 7: Latency-accuracy trade-off on SC-Base-1B.

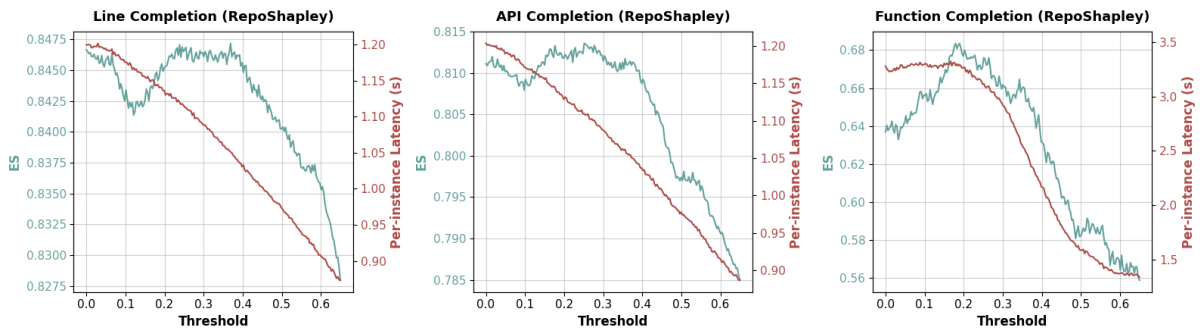


Figure 8: Latency-accuracy trade-off on SC-Base-3B.

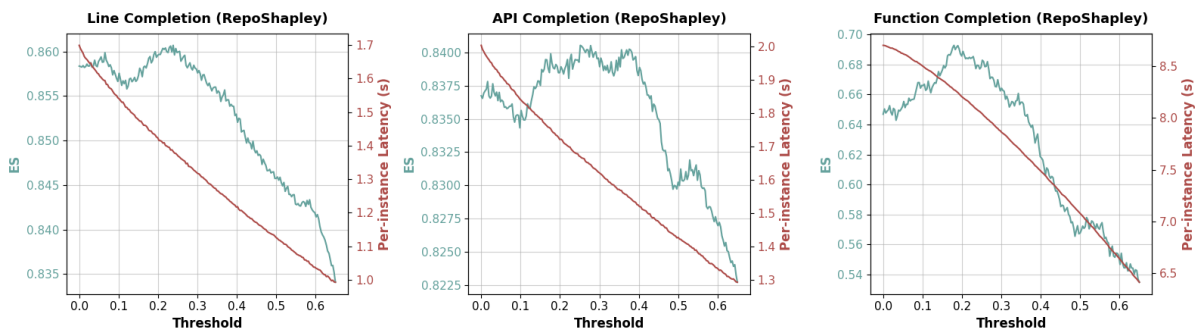


Figure 9: Latency-accuracy trade-off on SC-Base-7B.

<PFX>

```
import statistics
from datetime import datetime
class EventCollection(object):
    """Used to collect events found in logfiles..."""
    def __init__(self):
        self._events = {}
    def most_recent(self, items):
        return sorted(items, key=lambda e: e["end"],
                      reverse=True)[0]
    @property
    def complete_events(self):
        # ... (omitted for brevity if needed) ...
        return complete
    @property
    def incomplete_events(self):
        # ... (omitted for brevity) ...
        return incomplete
    def find_most_recent_start(self, event_id, end_ts):
        """ For a given event end marker, find the most recent
            start marker. """
        last = None
        for item in self._events[event_id].get("heads", []):
            start_ts = item["start"]
            if start_ts <= end_ts:
                if not last or start_ts > last["start"]:
                    last = item
        return last
    def add_event_end(self, event_id, end_ts):
        if event_id not in self._events:
            self._events[event_id] = {}
        if "tails" not in self._events[event_id]:
            self._events[event_id]["tails"] = [end_ts]
        else:
            self._events[event_id]["tails"].append(end_ts)
    def add_event_start(self, event_id, start_ts,
                       metadata=None,
                       metadata_key=None):
        # ... logic to add start markers ...
        pass
    def calculate_event_deltas(self):
        # ... logic to calc deltas ...
        pass
```

```
class SearchResultIndices(object):
    # ... index definitions ...
    pass
class LogEventStats(object):
    """Used to identify events within logs..."""
    def __init__(self, results, results_tag_prefix,
                 custom_idx=None):
        self.data = EventCollection()
        self.results = results
        self.results_tag_prefix = results_tag_prefix
        # ... init logic ...
    def run(self):
        """ Collect event start/end markers... """
        seq_idx = self.log_seq_idx
        end_tag = "{}-end".format(self.results_tag_prefix)
        for result in self.results.find_by_tag(end_tag):
            day = result.get(seq_idx.day)
            secs = result.get(seq_idx.secs)
            end = "{} {}".format(day, secs)
            end = datetime.strptime(end, "%Y-%m-%d
                                   %H:%M:%S.%f")
```

<SFX>

```
start = "{} {}".format(day, secs)
start = datetime.strptime(start, "%Y-%m-%d
                           %H:%M:%S.%f")
metadata = result.get(seq_idx.metadata)
meta_key = seq_idx.metadata_key
event_id = result.get(seq_idx.event_id)
self.data.add_event_start(event_id, start,
                          metadata=metadata,
                          metadata_key=meta_key)
self.data.calculate_event_deltas()
def get_top_n_events_sorted(self, max, reverse=True):
    # ... sorting logic ...
    return top_n_sorted
def get_event_stats(self):
    # ... stats logic ...
    return stats
```

<NEED>

Figure 10: **FIM instance.** The missing span inserts the `add_event_end` registration after parsing each end marker; the subsequent start-marker handling logic continues in the suffix.

Retrieved Candidates Pool & Selection Decisions

<C_1> [KEEP]

```
ts_formats = ["%Y-%m-%d %H:%M:%S.%f", "%Y-%m-%d %H:%M:%S"]
# ... (timestamp parsing logic) ...
def filter_by_age(cls, results, result_age_hours):
    # ...
    current = datetime.strptime(current, "%Y-%m-%d %H:%M:%S")
    # ...
```

<C_2> [DROP]

```
if message is not None:
    message = str(message).format(**fdict)
# ... message formatting utilities ...
@cached_property_attr
def type(self):
    """ Name of core.issues.IssueTypeBase object ... """
    # ...
```

<C_3> [DROP]

```
data_file = os.path.join(dtmp, 'data.txt')
# ... YAML generation for testing ...
class MyEventHandler(events.YEventCheckerBase):
    def __init__(self):
        super().__init__(yaml_defs_group='mygroup'...)
```

<C_4> [DROP]

```
def _override_keys(cls): return ['raises']
# ... YPropertyOverride logic ...
def message_with_format_dict_applied(self, property=None):
    # ...
```

<C_5> [DROP]

```
s = FileSearcher()
# ... Unit tests for SequenceSearchDef ...
def test_sequence_searcher_eof(self):
    # ...
```

<C_6> [DROP]

```
# ... Duplicate or similar unit tests ...
self.assertEqual(len(sections), 2)
# ...
```

<C_7> [DROP]

```
""" Returns a dict of ceph versions info ... """
out = self.cli_cache['ceph_versions']
# ... Ceph daemon logic ...
```

<C_8> [KEEP]

```
class YPropertySearchBase(YPropertyMappedOverrideBase):
    @classmethod
    def get_datetime_from_result(cls, result):
        """ This attempts to create a datetime object... """
        ts = result.get(1)
        # ...
        ts_formats = ["%Y-%m-%d %H:%M:%S.%f", ...]
        # ...
```

<C_9> [KEEP]

```
for r in sorted(_results, key=lambda i: i[0], reverse=True):
    # ... filtering logic ...
def apply_constraints(self, results):
    # ...
    results = self.filter_by_age(results, result_age_hours)
    # ...
```

<C_10> [DROP]

```
rename = 'unknown'
# ... package version check ...
def _get_bind_interfaces(self, type):
    """ For the given config network type ... """
    # ...
```

Figure 11: **10-chunk retrieved pool and selection.** REPOSHAPLEY identifies and keeps $\{c_1, c_8, c_9\}$ (timestamp utilities split across chunks) while dropping unrelated evidence. (Content abbreviated for display).

Ground Truth

```
self.data.add_event_end(
    result.get(seq_idx.event_id),
    end)
start_tag = "{}-start".format(
    self.results_tag_prefix)
for result in self.results.find_by_tag(
    start_tag):
    day = result.get(seq_idx.day)
    secs = result.get(seq_idx.secs)
```

REPOSHAPLEY (Ours)

```
self.data.add_event_end(
    result.get(seq_idx.event_id),
    end)
start_tag = "{}-start".format(
    self.results_tag_prefix)
for result in self.results.find_by_tag(
    start_tag):
    day = result.get(seq_idx.day)
    secs = result.get(seq_idx.secs)
```

Full-Retrieve Baseline

```
ts = "{} {}".format(day, secs)
ts_formats = ["%Y-%m-%d..."]
for format in ts_formats:
    try:
        end = datetime.strptime(ts,
            format)
        break
    except ValueError:
        continue
start_tag = "{}-start".format(
    self.results_tag_prefix)
for result in self.results.find_by_tag(
    start_tag):
    day = result.get(seq_idx.day)
    secs = result.get(seq_idx.secs)
```

Figure 12: **Completion comparison.** REPOSHAPLEY inserts the correct control-flow logic. Full-Retrieve is distracted and hallucinates redundant parsing logic.