#### A INCREMENTAL VERSION OF LIBRARIAN

We support sharing library components across clusters through the below incremental version of our algorithm. It processes clusters sequentially, rather than in parallel, and its output can depend on the ordering imposed upon the clusters.

$$\mathcal{L}_{0} = \varnothing$$

$$\{\rho'_{n}\} = \bigcup_{c \in \text{CLUSTER}(\{\rho_{n}\})} \{\rho'_{c,i}\}_{i=1}^{|c|}$$

$$(\Delta \mathcal{L}_{c}, \{\rho'_{c,i}\}_{i=1}^{|c|}) = \arg \min_{(\Delta \mathcal{L}, \{\rho'_{i}\}) \in \text{SAMPLE}_{k}(c; \mathcal{L}_{t-1})} \ell(\mathcal{L}_{t-1} \cup \Delta \mathcal{L}, \{\rho'_{i}\}_{i=1}^{|c|}) \quad \forall c \in \text{CLUSTER}(\{\rho_{n}\})$$

$$\mathcal{L}_{t} = \mathcal{L}_{t-1} \cup \Delta \mathcal{L}_{t}, \qquad \mathcal{L}^{\star} = \mathcal{L}_{|\text{CLUSTER}(\{\rho_{n}\})|}$$

## B ALGORITHM

594

595 596

597

598

600

605 606

607 608 609

610 611 612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631632633634

635 636

637

638

639

640

641

642

643

644

645

646

647

```
Algorithm 1 Refactoring Specialized Programs into a Joint Library
```

```
Require: Set of independent, specialized programs P_{initial} = \{\rho_1, \rho_2, \dots, \rho_n\}
Require: Sample Budget K
Ensure: Joint library \mathcal{L}_{final} and set of refactored programs P_{final}
 1: C \leftarrow \text{Cluster}(P_{initial})
 2: \mathcal{L}_{final} \leftarrow \emptyset, P_{final} \leftarrow \emptyset
 3: for all cluster c \in C do
                                                                                                                      ▶ Each cluster independently
 4:
             T_C \leftarrow \text{GroupIntoTuples}(c)
                                                                                                                       5:
             for all tuples \tau \in T_C do
                   \{f_{retrieved}\} \leftarrow \mathsf{RetrieveRelevantFromLibrary}(\mathcal{L}, \tau)
 6:
 7:
 8:
                   for i = 1 to K do
                                                                                                                                         \triangleright Sample k times
                          \begin{array}{l} (\{f_{new,i}\}, \{\rho_i'\}\}) \leftarrow \text{Sample}\left(f_{retrieved}, \tau\right) \\ S \leftarrow S \cup \{(\{f_{new,i}\}, \{\rho_i'\})\} \end{array} 
 9:
10:
             \begin{array}{l} (f_{best}, \{\rho_{best}'\}) \leftarrow \operatorname{RerankAndSelectBest}(S, \ell\left(\cdot\right)) \\ \mathcal{L}_{final} \cup \{f_{best}\} \\ P_{final} \cup \{\rho_{best}'\} \\ \text{end for} \end{array} 
11:
12:
                                                                                                                            ▶ Rerank using objective
13:
14:
15:
16: end for
17: return \mathcal{L}_{final}, P_{final}
```

# C EXPERIMENTAL SETUP

**Grouping Programs into Collections** To facilitate parallel application of LIBRARIAN and manage the dataset scale, we assume that semantically distant files will have minimal overlap in their optimal library functions. Therefore, our overall approach partitions the dataset into disjoint collections through clustering.

For **CodeContests**, these collections are constructed from an initial corpus of  $\sim$ 9k problems with Python solutions: We first filter these files, removing those whose selected canonical solution is under 10 lines (minimal refactoring potential). For the remaining 4596 solutions we use a language model to generate textual descriptions of canonical solutions—emphasizing reusable components—which are embedded using OpenAI's text-embedding-ada-002.

Agglomerative Clustering (Ward Jr, 1963) is subsequently applied to these embeddings to partition the files into a predefined number of initial clusters, in our case 120. To create uniformly sized experimental units, we subsample each such cluster to form collections of 30 files. This collection

 size was empirically chosen because it balanced between the runtime of LIBRARIAN without limiting compression. We select 10 collections that we then use to evaluate our methods.

For Transformers, since the number of models is on the lower end, we manually chose a set of popular LLM / VLM models and passed them to the agent in collections of 5 code sources.

**REGAL Baselines.** To evaluate the ability of our libraries to support reuse on new problems, we turn to the program synthesis tasks used in REGAL, where learned libraries are added to help the program synthesizer. We evaluate on the two domains published by the authors, Logo and Date. Because our clustering is inspired by REGAL but adds additional complexity, for fair comparison, we keep their setup the same and only augment the training using sample + MDL rerank procedure described in Section 5.

**Code Contests.** To evaluate LIBRARIAN on refactoring Code Contests we select 6 collections of 30 files (problems). In each collection we group the problems into tuples of size 3. We set the sample budget to be K=8, since our ablations show that with larger K we discover better libraries 2. We use the MDL objective for rankings. The model used for sampling is OpenAI's o4-mini (OpenAI, 2024). To obtain MDL scores we use Qwen 2.5 7B Instruct (Qwen et al., 2025) as a balance between quality, speed, and cost.

**Code Agents on Transformers and Diffusers Repositories.** To fairly evaluate performance on the task by state-of-the-art systems, we use coding agents that advertise long-context ability to reason about, write, and refactor code repositories. Specifically, we use Claude Code (Cl) (Anthropic, 2025) which uses the Opus 4.1.

We test whether code agents can refactor collections of code sources autonomously, without human intervention. Refactoring repositories with code agents involves planning and iterative (re-)implementation and testing. Code agents are prompted to perform each of these steps, with feedback from the unit tests. Agents must run and repair unit tests autonomously. We run coding agents multiple times per task, logging their progress in checklists stored in text files.

703

745746747748

749 750

751

752

754

755

The instruction provided to human evaluators is as follows:

```
704
705
           ## 1. Materials Provided
706
           You will be given a set of files for each example case:
708
           * **'original_programs.py'**: This file contains a set of 3 distinct Python programs, each presented with its
           corresponding problem description/query. This represents the "before" sta

* ***'v1.py'**: This file presents the first refactoring approach. It includes:

* The 3 refactored versions of the original programs.
710
                                section (e.g., 'codebank.py' or inline) containing helper functions. These helper functions
                might be retrieved from an existing common library or newly created during this refactoring.

* Either the retrieved or the new helper function sections may be _empty_, in case no programs existed in the codebank at the time or if no need helper functions were created by the LLM.
711
712
           * **'v2.py'**: This file presents the second, alternative refactoring approach. Similar to 'refactoring_v1.py
713
                  ', it includes:
                \star The 3 refactored versions of the original programs (using a different strategy than v1).
714
                * A "library" section with its own set of helper functions.
715
716
           **NOTE**: both refactorings had accuracy at least as good as the original programs.
717
           ## 2. Your Task
718
           Your primary task is to:
719
                **Review** the 'original_programs.py' to understand the initial code and the problems being solved.
**Analyze** both 'refactoring_v1.py' and 'refactoring_v2.py'. Pay close attention to how the original
720
                  programs have been restructured and what functionalities have been extracted into their respective
721
                  libraries
722
       25
                **Decide which refactoring (Version 1 or Version 2) you believe is "better," ** based on the evaluation
           3.
                  criteria provided below (or your own criteria!)
723
           ## 3. Evaluation Criteria: What to Consider for Your Choice
725
       29
           When comparing 'refactoring_v1.py' and 'refactoring_v2.py', please *consider* the following aspects to inform
                  your choice. The "better" refactoring should ideally excel in these areas:
726
           > Most importantly, make sure that the extracted functions are **actually reusable and not too specific.** If
727
                  the main programs are short, the refactoring is not immediately "better"! Try to think whether the extracted functions could actually be used in a different program down the line.
728
729
           * **Reusability of Helper Functions :**
                * **Generality:** Are the new helper functions general-purpose and potentially useful for *other,
730
                  different* programs and problems beyond the three presented?
                * **Reuse: ** How much were existing helper functions reused?
731
                  **Specificity:** Are the functions too specialized to the current set of problems, limiting their
732
                 broader applicability? _Avoid functions that are essentially just the original program broken out into a
                    "helper
733
                * Composability
           * **Maintainability:**
734
                * Readability & Understandability
                * Ease of Modification
735
                * Separation of Concerns
           ## 4. What NOT to Focus On:
737
           * **Comments:** Please disregard the presence or absence of comments in the code for this evaluation. These
738
                  are superficially generated by LLMs in some occasions and could be added manually after with a single
739
             **Minor Stylistic Differences:** Do not focus on trivial differences in variable naming or formatting,
740
                  unless they significantly impact readability or understanding
741
742
           For each example case, please provide:
743
               **Your Preferred Version: ** (e.g., "Version 1" or "Version 2")
744
```

Listing 1: Human Evaluation Instruction

### D ASYMPTOTIC BEHAVIOR RESULTS

Here we include the full graph of asymptotic behavior of all four scoring metrics (MDL, tokens, MI and CC), reporting the raw library metrics and the metric ratios. We can see that refactored programs have higher CC than the baseline and that optimizing CC as an objective does not decrese the other metrics. MDL performs best on total raw library usage as well as on average times a library function is used in the refactorings. MI ends up optimizing the number of library functions the best, but their reusability is below the average for the refactorings. Optimizing for tokens produces smaller libraries with less usage per function compared to optimizing MDL.

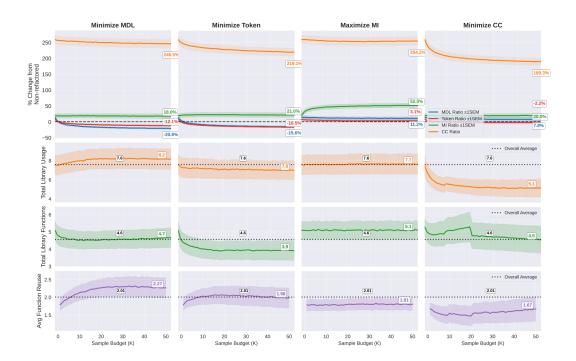


Figure 7: Asymptotic behavior of metrics for scoring libraries and refactorings (columns) varying refactoring budget (horizontal axes).

### E HUMAN STUDY DETAILS

We ran a user study with 12 participants where each participant had to judge 10 refactoring pairs. Each pair was comparing two out of three metrics, MDL, tokens, and MI. We showed the participants the original programs, as well as both of the refactoring versions (including the programs and the learned libraries). The participants were able to choose either version or say that the two refactorings are almost the same.

To quantify pairwise preferences between refactoring metrics, we employed the Bradley-Terry model, a standard framework for analyzing paired comparison data. We fit the model using maximum likelihood estimation with mutual information (MI) as the reference category ( $\pi_M I = 1$ ). To address potential noise in human judgments, we applied consensus-based filtering with a 75% threshold, retaining all responses for comparisons with low consensus (indicating genuine ambiguity) while excluding minority responses on high-consensus comparisons where the majority preference likely indicates the correct judgment. This conservative approach preserved 94.2% of responses while strengthening the statistical evidence for metric preferences, with MDL significantly outperforming MI (p = 0.67, 95% Confidence Interval: [0.56, 0.78]). Even without the filtering MDL preference over MI was statistically significant.

#### F BEST@K COMPRESSION IS A U-STATISTIC

We wish to estimate the expected compression ratio achieved by our sample + rerank method, which samples k candidate refactorings, discards any that do not pass the tests, and selects the one with the lowest score (total log-prob).

**Background on U-Statistics.** Let  $Z_1, \ldots, Z_n \overset{\text{i.i.d.}}{\sim} F$ . For a symmetric function  $h : \mathbb{Z}^k \to \mathbb{R}$ , the *U-statistic of order* k is defined as

$$U_n = \binom{n}{k}^{-1} \sum_{1 \le i_1 < \dots < i_k \le n} h(Z_{i_1}, \dots, Z_{i_k}). \tag{6}$$

By construction,

$$\mathbb{E}[U_n] = \mathbb{E}[h(Z_1, \dots, Z_k)],\tag{7}$$

so  $U_n$  is an unbiased estimator of the population quantity  $\theta = \mathbb{E}[h(Z_1, \dots, Z_k)]$ .

**Application to Best@k Compression.** Let each valid refactoring be a pair Z = (S, C), where S is the score and C is the compression ratio. Define the symmetric function

$$h_k(z_1, \dots, z_k) = C_{j^*}, \qquad j^* = \arg\min_{1 \le j \le k} S_j,$$
 (8)

the compression ratio of the lowest-score refactoring among k draws. The population target is then

$$\theta_k = \mathbb{E}[h_k(Z_1, \dots, Z_k)]. \tag{9}$$

Given n valid samples, our estimator is

$$\widehat{\theta}_k = \binom{n}{k}^{-1} \sum_{1 \le i_1 < \dots < i_k \le n} h_k(Z_{i_1}, \dots, Z_{i_k}). \tag{10}$$

**Proposition.**  $\widehat{\theta}_k$  is a U-statistic of order k with function  $h_k$ , and hence an unbiased estimator of  $\theta_k$ .

**Proof.** (1) Symmetry of the function  $h_k$ .  $h_k$  selects the compression associated with the lowest score among its k arguments. Permuting the inputs does not affect this outcome (ties can be resolved with a fixed, permutation-invariant rule). Thus  $h_k$  is symmetric.

(2) *U-statistic form.* By definition, a U-statistic of order k with kernel  $h_k$  is

$$U_n = \binom{n}{k}^{-1} \sum_{1 \le i_1 < \dots < i_k \le n} h_k(Z_{i_1}, \dots, Z_{i_k}),$$

which matches  $\widehat{\theta}_k$  exactly.

Therefore,  $\widehat{\theta}_k$  is a U-statistic of order k. By the unbiasedness property of U-statistics,

$$\mathbb{E}[\widehat{\theta}_k] = \theta_k.$$

Thus, our reported *best@k compression curves* provide unbiased estimates of the expected performance of the *sample* + *rerank* method.

## G CLUSTERING ANALYSIS: CODECONTESTS

We analyze the coherence of the clusters underlying collections in MINICODE-CodeContests. In particular, we compare clustering based on o4-mini generated file descriptions against task descriptions. Since task descriptions in competition coding problems are designed to hide the algorithmic approach needed to solve problem, we expect that clusters based on file descriptions are more coherent. We use Normalized Tag Instance Entropy and Herfindahl-Hirschman Index to evaluate clusterings. Figure 8 shows our clustering approach yields more thematically coherent clusters, evidenced by achieving lower entropy and higher HHI values across the entire tested range of N. We provide definitions of our measures below.

#### G.1 COLLECTION COHERENCE MEASURES

We use two measures to evaluate the thematic coherence of collections: Good collections should group files with a (1) concentrated and (2) identifiable set of shared *conceptual tags*, which for CodeContests are provided as ground truth (trees, graphs, etc.).

We provide the full definitions of the collection coherence measures below.

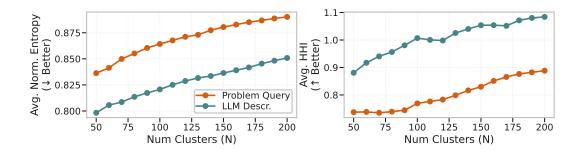


Figure 8: Clustering analysis of 4,596 Code Contest problems, comparing the thematic coherence of clusters formed using our proposed method versus REGAL-style clustering.

**Normalized Tag Instance Entropy:** This measures the concentration of tag *instances* within a collection C. Let  $p_i$  be the proportion of the i-th unique tag type among all tag instances in C, and  $D_C$  be the number of distinct tag types in C. If  $D_C > 1$ , the normalized entropy  $H_N$  is defined as:

$$H_N = -\frac{\sum_{i=1}^{D_C} p_i \log_2 p_i}{\log_2 D_C} \tag{11}$$

If  $D_C \le 1$ , then  $H_N = 0$ . Lower  $H_N$  (closer to 0) indicates higher thematic purity, meaning fewer tag types dominate the bulk of tag mentions.

**Herfindahl-Hirschman Index (HHI) for Problem Presence:** This measures tag concentration across distinct *problems* in a cluster C. Let  $s_t$  be the proportion of problems in C that include tag t (a problem contributes to  $s_t$  if t is one of its unique tags). A higher HHI signifies that the problems are collectively characterized by a smaller, more focused set of tags.

$$HHI = \sum_{t \in Tags(C)} s_t^2 \tag{12}$$

where Tags(C) represents the set of unique tags present in cluster C.

#### H BENCHMARK COMPARISON

We compare our benchmark, MINICODE, to similar benchmarks in Table 4. We define creativity and design as the need to explore diverse solutions in order to find the best solution possible. For example, optimizing for program correctness alone does not require exploring a large solutions space, whereas optimizing a program for speed would. In the case of compressing large files, we must explore the large space of shared abstractions afforded by libraries in order to maximize compression.

Table 4: Comparison of Code Benchmarks

Benchmark	Creativity/Design	Scale
SWE-bench (Jimenez et al., 2024)	Low	Repository
Commit-0 (Zhao et al., 2025)	Medium	Repository
RefactorBench (Gautam et al., 2025)	Low	File
ECCO (Waghjale et al., 2024)	High	Function
KernelBench (Ouyang et al., 2025)	High	Function
MINICODE(Ours)	High	Repository

## I FULL MINICODE CODECONTESTS RESULTS

We present the full agent scores for the CodeContests split in Table 5. The results are given both for each cluster of code sources, as well as averaged across clusters.

Cluster	Agent	Tokens	CC	Pass %	MDL	MDL %
0	original	9088	95	80.3	11745.85	100.0
	sonnet 3.7	18114	176	87.0	15005.18	127.7
	sonnet 4	11121	138	80.3	9901.53	84.3
	codex-mini	9321	95	80.3	9990.74	85.1
1	original	12531	255	89.7	13431.86	100.0
	sonnet 3.7	10470	239	96.7	8933.65	66.5
	sonnet 4	11325	298	96.7	8214.42	61.2
	codex-mini	12762	255	89.7	11798.73	87.8
2	original	14087	376	89.0	15012.77	100.0
	sonnet 3.7	17345	429	91.3	13145.02	87.6
	sonnet 4	14270	356	93.0	10522.66	70.1
	codex-mini	14318	376	89.0	13273.81	88.4
3	original	14261	246	90.3	13348.82	100.0
	sonnet 3.7	20749	241	97.7	15859.02	118.8
	sonnet 4	13433	197	80.7	11937.04	89.4
	codex-mini	14495	246	90.3	11616.41	87.0
4	original	17693	336	80.7	14665.16	100.0
	sonnet 3.7	29860	358	100.0	20666.52	141.0
	sonnet 4	18684	352	82.0	12801.21	87.3
	codex-mini	17923	336	80.7	12902.09	88.0
5	original	12588	286	92.0	12790.11	100.0
	sonnet 3.7	10580	128	99.3	8435.12	65.9
	sonnet 4	10416	155	99.3	9167.85	71.7
	codex-mini	12819	286	92.0	11086.19	86.7
6	original	11020	131	54.3	13540.41	100.0
	sonnet 3.7	21747	502	88.0	19446.07	143.6
	sonnet 4	11177	143	57.3	10492.00	77.5
	codex-mini	11251	131	54.3	11651.65	86.1
7	original	12301	180	80.0	12393.73	100.0
	sonnet 3.7	16390	166	91.0	13371.59	107.9
	sonnet 4	11625	150	85.7	9304.25	75.1
	codex-mini	12534	180	80.0	10549.04	85.1
Avg	original	12946	238	82.0	13366.09	100.0
	sonnet 3.7	18157	280	93.9	14357.77	107.4
	sonnet 4	12756	224	84.4	10292.62	77.1
	codex-mini	13178	238	82.0	11608.58	86.8

Table 5: Comparison of the pass rate and compression metrics of the original files, Claude Sonnet 4 and codex-mini refactorings across CodeContests clusters.

# J REFACTORING EXAMPLES OF LIBRARIAN ON CODE CONTESTS

### J.1 EXAMPLE 1

In code snippets 3, 2, 5, 4 one example of 2 refactoring versions. Specifically, the versions are both passing at least as many test cases as the original and they have the biggest difference in MDL among all the sample refactorings for that tuple. Sample + rerank filtering selected refactoring V2. You can observe that refactoring V1 introduces some problem specific functions like build\_max\_beauty\_perm(), while refactoring V2 sticks to more generally useful functions.

```
972
            # ==== NEW HELPER FUNCTIONS ====
973
           def compute_full_mask(i):
    """Return mask of all 1s of the bit-length of i."""
    return (1 << i.bit_length()) - 1</pre>
974
975
            def build_max_beauty_perm(n):
976
                 """Build permutation of 0..n maximizing sum of i^p[i].""" ans = [0] * (n + 1)
977
                 used = set()
                 for i in range(n, -1, -1): if i in used:
978
979
                           continue
                      mask = compute_full_mask(i)
j = i ^ mask
980
                      ans[i], ans[j] = j, i
981
                      used.add(i)
                 used.add(j)
beauty = sum(i ^ ans[i] for i in range(n + 1))
return ans, beauty
982
983
984
            {\tt def solve\_xor\_sum(u, v):}
985
                 Find shortest array whose xor is \boldsymbol{u} and sum is \boldsymbol{v} . Return list or None if impossible.
986
987
                 if u > v or (v - u) \% 2:
                 return None if u == v:
988
                 return [] if u == 0 else [u] x = (v - u) // 2
989
990
                 # try two elements
                 if ((u + x) ^ x) == u:
return [u + x, x]
991
                 # fallback to three elements
992
                 return [u, x, x]
993
            def build_trie(keys):
994
                 Build a binary trie with counts for 30-bit numbers.
995
        40
                 Each node: [left_index, right_index, count].
996
        41
                 tree = [[0, 0, 0]]
997
        43
                 for x in keys:
                     now = 0
998
                      tree[now][2] += 1
                      for i in range(29, -1, -1):
b = (x >> i) & 1
        46
999
                           if tree[now][b] == 0:
    tree[now][b] = len(tree)
1000
        48
1001
                                tree.append([0, 0, 0])
                           now = tree[now][b]
1002
                           tree[now][2] += 1
                 return tree
1003
1004
        55
            def trie_pop_min_xor(tree, x):
1005
                 Pop one key from trie to minimize x^key and return that minimal xor.
                 Decrements counts along the path.
1006
1007
        60
                 now = 0
                 res = 0
1008
                 for i in range(29, -1, -1):
b = (x >> i) & 1
nxt = tree[now][b]
1009
1010
                      if nxt and tree[nxt][2] > 0:
                          now = nxt
1011
                      else:
                      now = tree[now][b ^ 1]
res |= (1 << i)
tree[now][2] -= 1
1012
1013
                 return res
1014
```

Listing 2: Version 1, New Helpers

```
1026
1027
          # ######## PROGRAM: node 16:cc python 16 #########
1028
           from codebank import *
1029
           def main():
1030
               import sys
data = sys.stdin.readline()
1031
               if not data:
               return
n = int(data)
1032
1033
               perm, beauty = build_max_beauty_perm(n)
               print(beauty)
1034
               print(*perm)
1035
          if __name__ == "__main__":
1036
               main()
1037
          # ######## PROGRAM: node_19:cc_python_19 #########
1038
       20
          from codebank import *
1039
          def main():
1040
               import sys
data = sys.stdin.readline
1041
               A = Int(data())
A = list(map(int, data().split()))
P = list(map(int, data().split()))
1042
               trie = build_trie(P)
1043
               0 = [trie_pop_min_xor(trie, a) for a in A]
1044
               print(*0)
1045
          if __name__ == "__main__":
               main()
1046
          # ######## PROGRAM: node_25:cc_python_25 #########
1047
       37
38
1048
          from codebank import *
1049
               import sys
u, v = map(int, sys.stdin.readline().split())
res = solve_xor_sum(u, v)
1050
       41
1051
               if res is None:
1052
                   print(-1)
               else:
1053
                   print(len(res))
1054
                       print(*res)
1055
          if __name__ == "__main__":
1056
               main()
1057
```

Listing 3: Version 1, Refactored Programs

```
1062
            # ==== NEW HELPER FUNCTIONS ====
1063
            def compute_complement(i):
    return i ^ ((1 << i.bit_length()) - 1)</pre>
1064
            def trie_add(trie, x, max_bit):
    trie[0][2] += 1
1065
                 now = 0
for i in range(max_bit, -1, -1):
    bit = (x >> i) & 1
1066
1067
                      if trie[now][bit] == 0:
    trie[now][bit] = len(trie)
1068
                           trie.append([0, 0, 0])
1069
                      now = trie[now][bit]
trie[now][2] += 1
1070
1071
            def trie_find_min_xor(trie, x, max_bit):
        16
                 now = 0
1072
        18
                 ans = 0
                 for i in range(max_bit, -1, -1):
bit = (x >> i) & 1
1073
1074
                      if trie[now][bit] and trie[trie[now][bit]][2] > 0:
                           now = trie[now][bit]
1075
                           now = trie[now][bit ^ 1]
1076
                      ans |= (1 << i)
trie[now][2] -= 1
1077
                 return ans
1078
```

Listing 4: Version 2, New Helpers

```
1080
           # ######## PROGRAM: node_16:cc_python_16 ########
1081
           from codebank import *
1082
1083
           def main():
               imput = sys.stdin.readline
n = int(input())
ans = [-1] * (n + 1)
for i in range(n, -1, -1):
    if ans[i] == -1:
1084
1085
1086
                       z = compute_complement(i)
1087
                       ans[i] = z
ans[z] = i
1088
               m = sum(i ^ ans[i] for i in range(n + 1))
1089
               print(m)
1090
               print(*ans)
1091
       19
           if __name__ == "__main__":
1092
           # ######## PROGRAM: node_19:cc_python_19 ########
1093
1094
       24
           from codebank import *
1095
           def main():
               import sys
input = sys.stdin.readline
1096
1097
               n = int(input())
               A = list(map(int, input().split()))
P = list(map(int, input().split()))
1098
               1099
1100
               trie_add(trie, x, max_bit)
res = [trie_find_min_xor(trie, x, max_bit) for x in A]
1101
               print(*res)
1102
           if __name__ == "__main__":
1103
               main()
1104
           # ######## PROGRAM: node_25:cc_python_25 #########
1105
           from codebank import *
1106
          def main():
1107
               u, v = map(int, input().split())
if u > v or ((v - u) & 1):
    print(-1)
elif u == 0 and v == 0:
1108
1109
               print(0)
elif u == v:
1110
                    print(1)
1111
                    print(u)
1112
               else:
                     u = (v - u) // 2
1113
                    if (w & u) == 0:
d = u + w
1114
                        print(2)
1115
       60
                         print(d, w)
                    else:
1116
                        print(3)
                         print(u, w, w)
1117
1118
           if __name__ == "__main__":
               main()
1119
```

Listing 5: Version 2, Refactored Programs

# J.2 EXAMPLE 2

1127 1128

11301131

1132

1133

In code snippets 7, 6, 9, 8 is another example of 2 refactorings where V1 was better according to LIBRARIAN. We can observe that V2 creates helper functions that are overly specific to the problem. You can see that refactoring V2 introduces overly specialized functions like dijkstra\_special() or compute\_min\_moves\_opposite\_parity(). In comparison, refactoring V1 generates only general versions of these functions (e.g. dijkstra()).

```
1134
1135
             # ==== NEW HELPER FUNCTIONS ====
             def read_ints():
1136
                   return list(map(int, input().split()))
1137
             def build_adj_undirected(n, edges):
1138
                   adj = [[] for _ in range(n)]
for u, v, w in edges:
    adj[u].append((v, w))
1139
                   adj[v].append((u, w))
return adj
1140
1141
             def dijkstra(adj, src):
    from heapq import heappush, heappop
    INF = 10**18
    n = len(adj)
1142
1143
                   dist = [INF]*n
1144
                   parent = [-1]*n
dist[src] = 0
1145
         19
                   heap = [(0, src)]
while heap:
d, u = heappop(heap)
         20
21
1146
1147
                        if d > dist[u]:
                        continue
for v, w in adj[u]:
nd = d + w
if nd < dist[v]:
dist[v] = nd
parent[v] = u
1148
1149
1150
1151
                                   heappush(heap, (nd, v))
1152
                   return dist, parent
1153
             def reconstruct_path(parent, dest):
                   path = []
1154
                   u = dest
1155
                   while u != -1:
                        path.append(u+1)
1156
                   u = parent[u]
return path[::-1]
1157
             def multi_source_bfs(neighbors, sources):
1158
        41
                   from collections import deque
1159
                   n = len(neighbors)
dist = [-1]*n
dq = deque()
1160
                   for u in sources:
    if dist[u] == -1:
        dist[u] = 0
1161
1162
         48
                             dq.append(u)
1163
                   while dq:
                        u = dq.popleft()
1164
                        for v in neighbors[u]:
    if dist[v] == -1:
        dist[v] = dist[u] + 1
1165
1166
                                   dq.append(v)
                   return dist
1167
```

Listing 6: Version 1, New Helpers

```
1188
                          # ######## PROGRAM: node_16:cc_python_16 ########
1189
                          from codebank import *
1190
                          def main()
1191
                                   Import neapy, n, m = read_ints() edges = [(u-1, v-1, w) for u, v, w in (read_ints() for _ in range(m))] adj = build_adj_undirected(n, edges)
INF = 10**20
dist = [INF]*n
1192
1193
1194
                                    dist[0] = 0
last_w = [0]*n
heap = [(0, 0)]
1195
1196
                                    while heap:
d, u = heapq.heappop(heap)
1197
                                             if d > dist[u]:
1198
                                             continue
# record last edges
1199
                  19
                                            for v, w in adj[u]:
1200
                                              last_w[v] = w
# expand two-edge moves
1201
                                               for v, w1 in adj[u]:
tw = last_w[v]
                                                       tw = last_w[v]
for x, w2 in adj[v]:
    nd = d + (tw + w2)**2
    if nd < dist[x]:
        dist[x] = nd</pre>
1202
1203
1204
1205
                 29
                                                                              heapq.heappush(heap, (nd, x))
                                     out = []
1206
                                    for x in dist:
                                    out.append(str(x if x < INF else -1)) print("".join(out))
1207
1208
                         if __name__ == "__main__":
    main()
1209
                 37
38
1210
                          # ######## PROGRAM: node_17:cc_python_17 #########
1211
                          from codebank import *
1212 41
                          def main():
1213
                43
                                   n, m = read_ints()
                                    n, m = read_ints()
edges = [(u-1, v-1, w) for u, v, w in (read_ints() for _ in range(m))]
adj = build_adj_undirected(n, edges)
1214
                                   dist, parent = dijkstra(adj, 0)
if dist[n-1] >= 10**18:
1215
1216
                 48
                                             print(-1)
                                    else:
1217
                                              path = reconstruct_path(parent, n-1)
                                               print(*path)
1218
1219
                          if __name__ == "__main__":
                                   main()
1220
                 55
                          # ######## PROGRAM: node_19:cc_python_19 #########
1221
                        from codebank import *
1222
1223
                 60
                          def main():
                                   n = int(input())
1224
                                    a = read_ints()
                                    # build reversed graph: for each move i->j, add edge j->i
1225
                                    neighbors = [[] for _ in range(n)]
                                    for i, val in enumerate(a):
    for j in (i - val, i + val):
        if 0 <= j < n:</pre>
1226
1227
                                                                 neighbors[j].append(i)
                                   neighbors[j].append(1)

# BFS from all even and all odd positions separately
even_sources = [i for i, val in enumerate(a) if val % 2 == 0]
odd_sources = [i for i, val in enumerate(a) if val % 2 == 1]
dist_even = multi_source_bfs(neighbors, even_sources)
dist_odd = multi_source_bfs(neighbors, odd_sources)
# froudd = [i] answer is dist to nearest even => dist even; of the content 
1228
1229
1230
1231
                                    # for odd a[i], answer is dist to nearest even => dist_even; else dist_odd
ans = [dist_even[i] if a[i] % 2 == 1 else dist_odd[i] for i in range(n)]
1232
                                     print(*ans)
1233
                          if __name__ == "__main__":
1234
                                    main()
```

Listing 7: Version 1, Refactored Programs

```
1242
1243
            # ==== NEW HELPER FUNCTIONS ====
            def read_ints():
1244
                 return list(map(int, input().split()))
1245
            def build_undirected_weighted_graph(n, m):
1246
                 from collections import defaultdict
adj = defaultdict(list)
1247
                 for _ in range(m):
                     u, v, w = read_ints()
u -= 1; v -= 1
1248
1249
                      adj[u].append((v, w))
                      adj[v].append((u, w))
1250
                 return adj
1251
            def diikstra(adi. src. n):
                 import heapq
1252
                 INF = 10**18
dist = [INF]*n
1253
                 parent = [-1]*n
visited = [False]*n
dist[src] = 0
1254
1255
                 heap = [(0, src)]
while heap:
    d, u = heapq.heappop(heap)
1256
1257
                      if visited[u]:
                      continue
visited[u] = True
1258
                      for v, w in adj.get(u, ()):
    nd = d + w
    if nd < dist[v]:</pre>
1259
1260
                                dist[v] = nd
parent[v] = u
1261
                                heapq.heappush(heap, (nd, v))
1262
                 return dist, parent
1263
            def reconstruct_path(parent, dest):
1264
                 path = []
while dest != -1:
1265
        40
                      path.append(dest+1)
                 dest = parent[dest]
return path[::-1]
1266
        41
1267
        43
            def dijkstra_special(e, n, src):
1268
                 import heapq
INF = 10**18
d = [INF]*n
1269
1270
        48
                 d[src] = 0
                 heap = [(0, src)]
                 1271
1272
1273
                      td = {}
for u, w in e.get(v, ()):
    td[u] = w
1274
1275
                      for u, w1 in td.items():
	for x, w2 in e.get(u, ()):
		cost = cd + (w1 + w2)**2
1276
1277
                                if cost < d[x]:
    d[x] = cost</pre>
        60
1278
                                      heapq.heappush(heap, (cost, x))
        63
                 return d
1279
1280
        65
            def compute_min_moves_opposite_parity(a):
                 from collections import deque
1281
                 n = len(a)
                 go = [[] for _ in range(n)]
ans = [-1]*n
1282
        69
                 q = deque()
1283
                 for i, val in enumerate(a):
    for j in (i - val, i + val):
        if 0 <= j < n:</pre>
1284
1285
                                if (a[j] % 2) != (val % 2):
ans[i] = 1
1286
                                     q.append(i)
1287
                                      break
                                else:
1288
        70
                                      go[j].append(i)
                 while q:
    u = q.popleft()
1289
                      for v in go[u]:
    if ans[v] == -1:
1290
1291
                                ans[v] = ans[u] + 1
                                q.append(v)
1292
                 return ans
1293
```

Listing 8: Version 2, New Helpers

```
1296
1297
           # ######## PROGRAM: node_16:cc_python_16 #########
1298
           from codebank import *
1299
           def main():
1300
               n, m = read_ints()
e = {}
for _ in range(m):
1301
                   u, v, w = read_ints()
u -= 1; v -= 1
1302
1303
                    e.setdefault(u, \ []).append((v, \ w))
               e.setdefault(v, []).append((u, w))
d = dijkstra_special(e, n, 0)
1304
               print(" ".join(str(-1 if x \ge 10**18 else int(x)) for x in d))
1305
1306
           if __name__ == "__main__":
               main()
1307
           # ######## PROGRAM: node_17:cc_python_17 #########
1308
           from codebank import *
1309
1310
       24
           def main():
               n, m = read_ints()
adj = build_undirected_weighted_graph(n, m)
1311
               dist, parent = dijkstra(adj, 0, n)
if dist[n-1] >= 10**18:
1312
1313
                   print(-1)
               else:
1314
                    path = reconstruct_path(parent, n-1)
print(" ".join(map(str, path)))
1315
1316
           if __name__ == "__main__":
               main()
1317
           # ######## PROGRAM: node_19:cc_python_19 #########
1318
           from codebank import *
1319
1320
       41
           def main():
               n = int(input())
1321
               a = read_ints()
               ans = compute_min_moves_opposite_parity(a)
print(" ".join(map(str, ans)))
1322
1323
           1324
1325
```

Listing 9: Version 2, Refactored Programs

#### K OBFUSCATION EXAMPLE

1326 1327 1328

1329 1330 1331

1332 1333

1334

We provide the full source code, both refactored and obfuscated, for the MDL and token comparison here.

Listing 10: Refactored code from modeling\_llama.py from the refactored Transformers repository. from typing import Optional, Union

```
1335
           import torch
1336
           from torch import nn
1337
           from ...cache_utils import Cache, DynamicCache
1338
           from ...generation import GenerationMixin
1339
           from ...masking_utils import create_causal_mask
           from ...modeling_layers import (
1340
               {\tt GenericForQuestionAnswering,}
1341
               GenericForSequenceClassification,
1342
               GenericForTokenClassification,
1343
           from ...modeling_outputs import (
1344
               BaseModelOutputWithPast,
               CausalLMOutputWithPast,
1345
           from ...modeling_utils import PreTrainedModel
1346
           \textcolor{red}{\textbf{from}} \ \dots \\ \textbf{processing\_utils} \ \textcolor{red}{\textbf{import}} \ \texttt{Unpack}
1347
           from ...utils import TransformersKwargs, auto_docstring, can_return_tuple, logging
           from ...utils.generic import check_model_inputs
1348
           from .configuration_llama import LlamaConfig
1349
           from ..shared_library import (
```

```
1350
             rotate_half,
1351
              apply_rotary_pos_emb,
              repeat_kv,
1352
              eager_attention_forward,
1353
              BaseMLP,
1354
              BaseRotaryEmbedding,
1355
             BaseAttention,
             BaseDecoderLayer,
1356
         )
1357
1358
          logger = logging.get_logger(__name__)
1359
1360
          class LlamaRMSNorm(RMSNorm):
1361
             pass
1362
1363
          class LlamaRotaryEmbedding(BaseRotaryEmbedding):
1364
             pass
1365
1366
          class LlamaMLP(BaseMLP):
              def __init__(self, config):
1367
                  super().__init__(config, mlp_bias=config.mlp_bias)
1368
1369
          class LlamaAttention(BaseAttention):
              def __init__(self, config: LlamaConfig, layer_idx: int):
1370
                  super().__init__(
1371
                      config=config,
                      layer_idx=layer_idx,
1372
                      attention_bias=config.attention_bias,
1373
                      sliding_window=None
                 )
1374
1375
          class LlamaDecoderLayer(BaseDecoderLayer):
1376
              def __init__(self, config: LlamaConfig, layer_idx: int):
1377
                  super().__init__(
                      config=config,
1378
                      layer_idx=layer_idx,
1379
                      norm_class=LlamaRMSNorm,
                      mlp_class=LlamaMLP,
1380
                      attention\_class = Llama Attention
1381
1382
1383
          @auto_docstring
          class LlamaPreTrainedModel(PreTrainedModel):
1384
              config: LlamaConfig
1385
              base_model_prefix = "model"
              supports_gradient_checkpointing = True
1386
              _no_split_modules = ["LlamaDecoderLayer"]
1387
              _skip_keys_device_placement = ["past_key_values"]
             _supports_flash_attn = True
1388
              _supports_sdpa = True
1389
             _supports_flex_attn = True
1390
              _can_compile_fullgraph = True
1391
              _supports_attention_backend = True
              _can_record_outputs = {
1392
                  "hidden_states": LlamaDecoderLayer,
1393
                  "attentions": LlamaAttention,
             }
1394
1395
          @auto docstring
1396
          class LlamaModel(LlamaPreTrainedModel):
1397
              def __init__(self, config: LlamaConfig):
                  super().__init__(config)
self.padding_idx = config.pad_token_id
1398
1399
                  self.vocab_size = config.vocab_size
                  self.embed_tokens = nn.Embedding(config.vocab_size, config.hidden_size, self.padding_idx)
1401
                  self.layers = nn.ModuleList(
                      [LlamaDecoderLayer(config, layer_idx) for layer_idx in range(config.num_hidden_layers)]
1402
1403
                  self.norm = LlamaRMSNorm(config.hidden_size, eps=config.rms_norm_eps)
                  self.rotary_emb = LlamaRotaryEmbedding(config=config)
```

```
1404
                  self.gradient_checkpointing = False
1405
                  self.post_init()
1406
1407
              @check_model_inputs
              @auto_docstring
1408
              def forward(
1409
                  self,
                  input_ids: Optional[torch.LongTensor] = None,
1410
                  attention_mask: Optional[torch.Tensor] = None,
1411
                  position_ids: Optional[torch.LongTensor] = None,
                  past_key_values: Optional[Cache] = None,
1412
                  inputs_embeds: Optional[torch.FloatTensor] = None,
1413
                  cache_position: Optional[torch.LongTensor] = None,
                  use_cache: Optional[bool] = None,
1414
                  **kwargs: Unpack[TransformersKwargs],
1415
              ) -> BaseModelOutputWithPast:
                  if (input_ids is None) ^ (inputs_embeds is not None):
1416
                      raise ValueError("You must specify exactly one of input_ids or inputs_embeds")
1417
1418
                  if inputs embeds is None:
                      inputs embeds: torch.Tensor = self.embed tokens(input ids)
1419
                  if use_cache and past_key_values is None:
1420
                      past_key_values = DynamicCache(config=self.config)
1421
                  if cache_position is None:
1422
                      past_seen_tokens = past_key_values.get_seq_length() if past_key_values is not None else 0
1423
                      cache_position: torch.Tensor = torch.arange(
                          past_seen_tokens, past_seen_tokens + inputs_embeds.shape[1], device=inputs_embeds.device
1424
1425
                  if position_ids is None:
1426
                      position_ids = cache_position.unsqueeze(0)
1427
                  causal_mask = create_causal_mask(
1428
                      config=self.config,
1429
                      input_embeds=inputs_embeds,
                      attention_mask=attention_mask,
1430
                      cache_position=cache_position,
1431
                      past_key_values=past_key_values,
                      position_ids=position_ids,
1432
1433
                  hidden\_states = inputs\_embeds
1434
                  position_embeddings = self.rotary_emb(hidden_states, position_ids)
1435
                  for decoder_layer in self.layers[: self.config.num_hidden_layers]:
1436
                      hidden_states = decoder_layer(
1437
                          hidden_states,
                          attention_mask=causal_mask,
1438
                          position_ids=position_ids,
1439
                          past_key_values=past_key_values,
                          cache_position=cache_position,
1440
                          position_embeddings=position_embeddings,
1441
                          **kwargs,
1442
1443
                  hidden_states = self.norm(hidden_states)
                  return BaseModelOutputWithPast(
1444
                      last_hidden_state=hidden_states,
1445
                      past_key_values=past_key_values,
                  )
1446
1447
          @auto_docstring
1448
          class LlamaForCausalLM(LlamaPreTrainedModel, GenerationMixin):
1449
              _tied_weights_keys = ["lm_head.weight"]
_tp_plan = {"lm_head": "colwise_rep"}
1450
              _pp_plan = {"lm_head": (["hidden_states"], ["logits"])}
1451
              def __init__(self, config):
    super().__init__(config)
1452
1453
                  self.model = LlamaModel(config)
                  self.vocab_size = config.vocab_size
                  self.lm_head = nn.Linear(config.hidden_size, config.vocab_size, bias=False)
1455
                  self.post_init()
1456
1457
              def set decoder(self. decoder):
                  self.model = decoder
```

```
1458
1459
              def get_decoder(self):
                  return self.model
1460
1461
              @can_return_tuple
              @auto_docstring
1462
              def forward(
1463
                 self,
                  input_ids: Optional[torch.LongTensor] = None,
1464
                 attention_mask: Optional[torch.Tensor] = None,
1465
                 position_ids: Optional[torch.LongTensor] = None,
                 past_key_values: Optional[Cache] = None,
1466
                  inputs_embeds: Optional[torch.FloatTensor] = None,
1467
                 labels: Optional[torch.LongTensor] = None,
                 use cache: Optional[bool] = None.
1468
                 cache_position: Optional[torch.LongTensor] = None,
1469
                  logits_to_keep: Union[int, torch.Tensor] = 0,
                  **kwargs: Unpack[TransformersKwargs],
1470
             ) -> CausalLMOutputWithPast:
1471
1472
                 Example:
1473
                 '''python
                 >>> from transformers import AutoTokenizer, LlamaForCausalLM
1474
1475
                 >>> model = LlamaForCausalLM.from pretrained("meta-llama/Llama-2-7b-hf")
                 >>> tokenizer = AutoTokenizer.from_pretrained("meta-llama/Llama-2-7b-hf")
1476
1477
                 >>> prompt = "Hey, are you conscious? Can you talk to me?"
                 >>> inputs = tokenizer(prompt, return_tensors="pt")
1478
1479
                 >>> # Generate
                 >>> generate_ids = model.generate(inputs.input_ids, max_length=30)
1480
                 >>> tokenizer.batch_decode(generate_ids, skip_special_tokens=True, clean_up_tokenization_spaces=False)[0]
1481
                 "Hey, are you conscious? Can you talk to me?\nI'm not conscious, but I can talk to you."
1482
                 outputs: BaseModelOutputWithPast = self.model(
1483
                     input_ids=input_ids,
                     attention_mask=attention_mask,
1484
                     position_ids=position_ids,
1485
                     past_key_values=past_key_values,
                      inputs_embeds=inputs_embeds,
1486
                     use_cache=use_cache,
1487
                     cache_position=cache_position,
1488
                      **kwargs,
1489
                 hidden_states = outputs.last_hidden_state
1490
                  slice_indices = slice(-logits_to_keep, None) if isinstance(logits_to_keep, int) else logits_to_keep
1491
                  logits = self.lm_head(hidden_states[:, slice_indices, :])
1492
                 loss = None
1493
                     loss = self.loss_function(logits=logits, labels=labels, vocab_size=self.config.vocab_size, **kwargs)
1494
1495
                 return CausalLMOutputWithPast(
                     loss=loss,
1496
                      logits=logits,
1497
                     past_key_values=outputs.past_key_values,
                     hidden_states=outputs.hidden_states,
1498
                     attentions=outputs.attentions,
1499
1501
         class LlamaForSequenceClassification(GenericForSequenceClassification, LlamaPreTrainedModel): ...
1502
1503
         class LlamaForOuestionAnswering(GenericForOuestionAnswering, LlamaPreTrainedModel):
1504
             base model prefix = "transformer
1505
         class LlamaForTokenClassification(GenericForTokenClassification, LlamaPreTrainedModel): ...
1506
1507
          __all__ = [
1508
              "LlamaForCausalLM",
1509
              "LlamaModel",
              "LlamaPreTrainedModel"
1510
              "LlamaForSequenceClassification",
              "LlamaForQuestionAnswering",
              "LlamaForTokenClassification",
```

```
1512
                          ]
1513
1514
                          Listing 11: Obfuscated code from modeling_llama.py from the refactored Transformers repository.
1515
                          from typing import Optional as O, Union as U
1516
                          import torch as T
1517
                          from torch import nn as N
                          from ...cache_utils import Cache as C,DynamicCache as D
1518
                          from ...generation import GenerationMixin as G
1519
                          from ...masking_utils import create_causal_mask as m
                          1520
                                         GenericForTokenClassification as K
1521
                          from ...modeling_outputs import BaseModelOutputWithPast as B,CausalLMOutputWithPast as L
1522
                          from ...modeling_utils import PreTrainedModel as P
                          from ...processing_utils import Unpack as X
1523
                          from ...utils import TransformersKwargs as W,auto_docstring as a,can_return_tuple as r,logging
                          from ...utils.generic import check_model_inputs as c
1524
                          from .configuration_llama import LlamaConfig as F
1525
                          l = logging.get_logger(\__name\__)
1527
                          class Z(R):pass
                          class I(E):pass
1528
                          class J(M):
                            def __init__(s,g):super().__init__(g,mlp_bias=g.mlp_bias)
                          class H(A):
1530
                            def __init__(s,g:F,i:int):super().__init__(config=g,layer_idx=i,attention_bias=g.attention_bias,sliding_window=
                          class V(Y):
1532
                            \label{lem:config} \textbf{def} \ \_\texttt{init}\_(\texttt{s,g:F,i:int}): \textbf{super}(). \ \_\texttt{init}\_(\texttt{config=g,layer\_idx=i,norm\_class=Z,mlp\_class=J,attention\_class=H})
1533
                          @a
                          class o(P):
1534
                            config:F
1535
                            base_model_prefix="model"
                             {\tt supports\_gradient\_checkpointing=True}
1536
                            _no_split_modules=["LlamaDecoderLayer"]
1537
                            _skip_keys_device_placement=["past_key_values"]
                            _supports_flash_attn=True
1538
                             _supports_sdpa=True
1539
                            _supports_flex_attn=True
                            _can_compile_fullgraph=True
1540
                            \verb|_supports_attention_backend=True|
1541
                             _can_record_outputs={"hidden_states":V, "attentions":H}
1542
                          class u(o):
1543
                            def __init__(s,g:F):
                               super().__init__(g)
1544
                               s.padding_idx=g.pad_token_id
                               s.vocab_size=g.vocab_size
                               s.embed_tokens=N.Embedding(g.vocab_size,g.hidden_size,s.padding_idx)
                               s.layers=N. \\ ModuleList([V(g,i) \\ for i \\ in \\ range(g.num\_hidden\_layers)])
                               s.norm=Z(g.hidden_size,eps=g.rms_norm_eps)
                               s.rotary_emb=I(config=g)
1548
                               s.gradient checkpointing=False
1549
                              s.post_init()
1550
1551
                            def forward(s,input_ids:0[T.LongTensor]=None,attention_mask:0[T.Tensor]=None,position_ids:0[T.LongTensor]=None,
                                         \verb|past_kev_values: O[C]=None, inputs_embeds: O[T.FloatTensor]=None, cache\_position: O[T.LongTensor]=None, use\_cache: O[T.FloatTensor]=None, 
1552
                                         0[boo1]=None, **k: X[W])->B:
1553
                               if(input_ids is None)^(inputs_embeds is not None):raise ValueError("You must specify exactly one of input_ids
                                          or inputs embeds")
1554
                               if inputs_embeds is None:inputs_embeds:T.Tensor=s.embed_tokens(input_ids)
1555
                               if use_cache and past_key_values is None:past_key_values=D(config=s.config)
1556
                               if cache_position is None:
                                  p=past_key_values.get_seq_length()if past_key_values is not None else 0 \,
1557
                                  {\tt cache\_position:T.Tensor=T.arange(p,p+inputs\_embeds.shape[1],device=inputs\_embeds.device)}
                                \textbf{if} \ position\_ids \ \textbf{is} \ \textbf{None}: position\_ids=cache\_position.unsqueeze(0) 
1558
                               f=m(config=s.config,input\_embeds=inputs\_embeds,attention\_mask=attention\_mask,cache\_position=cache\_position,attention\_mask=attention\_mask,cache\_position=cache\_position,attention\_mask=attention\_mask,cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=cache\_position=ca
1559
                                          past_key_values=past_key_values,position_ids=position_ids)
                               h=inputs_embeds
                               e=s.rotary_emb(h,position_ids)
1561
                                for \ d \ in \ s.layers[:s.config.num\_hidden\_layers]: h=d(h, attention\_mask=f, position\_ids=position\_ids, past\_key\_values) ) and the property of the proper
                                         =past_key_values,cache_position=cache_position,position_embeddings=e,**k)
                               h=s.norm(h)
                               return B(last_hidden_state=h,past_key_values=past_key_values)
1564
                          class t(o,G):
                            _tied_weights_keys=["lm_head.weight"]
                             _tp_plan={"lm_head":"colwise_rep"}
```

```
1566
                                                               _pp_plan={"lm_head":(["hidden_states"],["logits"])}
1567
                                                               def __init__(s,g):
                                                                    super().__init__(g)
1568
                                                                    s.model=u(g)
1569
                                                                     s.vocab_size=g.vocab_size
                                                                    s.lm_head=N.Linear(g.hidden_size,g.vocab_size,bias=False)
1570
                                                                    s.post_init()
1571
                                                               def set_decoder(s,d):s.model=d
                                                               def get_decoder(s):return s.model
1572
1573
                                                               \textbf{def} \ \ forward(\texttt{s}, \texttt{input\_ids}: \texttt{O[T.LongTensor]} = \texttt{None}, \\ \textbf{attention\_mask}: \texttt{O[T.Tensor]} = \texttt{None}, \\ \textbf{position\_ids}: \texttt{O[T.LongTensor]} = \texttt{None}, \\ \textbf{optimize} = \texttt{None}, \\ \textbf{
1574
                                                                                          past_key_values:O[C]=None,inputs_embeds:O[T.FloatTensor]=None,labels:O[T.LongTensor]=None,use_cache:O[bool
1575
                                                                                          ]=None, cache_position:O[T.LongTensor]=None, logits_to_keep:U[int,T.Tensor]=0,**k:X[W])->L:
                                                                    o: B=s.model (input\_ids=input\_ids, attention\_mask=attention\_mask, position\_ids=position\_ids, past\_key\_values=attention\_mask, position\_ids=position\_ids, past\_key\_values=attention\_ids, 
1576
                                                                                         \verb|past_key_values, inputs_embeds=inputs_embeds, use_cache=use_cache, cache_position=cache\_position, **k)|
1577
                                                                     h=o.last\_hidden\_state
                                                                     i=slice(-logits_to_keep,None)if isinstance(logits_to_keep,int)else logits_to_keep
1578
                                                                    g=s.lm_head(h[:,i,:])
1579
                                                                      \textbf{if labels is not None}: n=s.loss\_function(logits=g,labels=labels,vocab\_size=s.config.vocab\_size,**k) \\ 
1580
                                                                     \textbf{return L} (loss=n, logits=g, past\_key\_values=o.past\_key\_values, hidden\_states=o.hidden\_states, attentions=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key\_values=o.past\_key
1581
                                                                                         attentions)
                                                          class b(S,o):...
1582
                                                          class x(Q,o):base_model_prefix="transformer"
1583
                                                        class y(K,o):...
__all__=["t","u","o","b","x","y"]
1584
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
```