

555 **Checklist**

- 556 1. For all authors...
 - 557 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
558 contributions and scope? [Yes] Please see Sections 2, 3, and 4.
 - 559 (b) Did you describe the limitations of your work? [Yes] Please see the paragraph on
560 **Limitations** in Section 2.
 - 561 (c) Did you discuss any potential negative societal impacts of your work? [N/A] The
562 CodeNet dataset is about code written for pedagogical purposes. We believe that
563 it does not have any negative societal impact. On the contrary, we are launching
564 a challenge/contest based on the CodeNet dataset with the Global Women in Data
565 Science organization with presence in over 50 countries to promote diversity, inclusion
566 and data science education in the field of AI for Code.
 - 567 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
568 them? [Yes] The paper presents a dataset with code samples submitted by students to
569 simple programming problems. The dataset neither does harm to living beings, nor
570 raise any security and economic concerns, human rights and surveillance issues, nor
571 damage the environment, nor deceive people and damage their livelihood. We have
572 anonymized each submitter's user id, and tried filtering offensive words. We have also
573 followed the term of service of the website from which we download the dataset.
- 574 2. If you are including theoretical results...
 - 575 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - 576 (b) Did you include complete proofs of all theoretical results? [N/A]
- 577 3. If you ran experiments (e.g. for benchmarks)...
 - 578 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
579 imental results (either in the supplemental material or as a URL)? [Yes] The source
580 code and instructions of the experiments are available in the model-experiments folder
581 at https://github.com/IBM/Project_CodeNet, when third-party licenses allow.
582 The datasets used in the experiments are available in <https://developer.ibm.com/technologies/artificial-intelligence/data/project-codenet>.
 - 583 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
584 were chosen)? [Yes] Please see Section 8 and appendix D, appendix E, and appendix F
585 in the supplementary materials.
 - 586 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
587 ments multiple times)? [Yes] Please see Section 8.
 - 588 (d) Did you include the total amount of compute and the type of resources used (e.g.,
589 type of GPUs, internal cluster, or cloud provider)? [Yes] Please see Section 8 and
590 appendix D, appendix E, and appendix F in the supplementary materials.
- 592 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 593 (a) If your work uses existing assets, did you cite the creators? [N/A] Our work is
594 creating/releasing new assets
 - 595 (b) Did you mention the license of the assets? [Yes] The license of the dataset is
596 CDLA Permissive v2.0. It is mentioned in <https://developer.ibm.com/technologies/artificial-intelligence/data/project-codenet> and
597 <https://www.linuxfoundation.org/press-release/enabling-easier-collaboration-on-open-data-for-ai-and-ml-with-cdla-permissive-20/>.
 - 601 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
602 The new assets are available at <https://developer.ibm.com/technologies/artificial-intelligence/data/project-codenet>.
 - 604 (d) Did you discuss whether and how consent was obtained from people whose data you're
605 using/curating? [Yes] We have looked into the terms of service, and ensured that the
606 code samples can be used for research purposes and we also contacted the respective
607 communities.

608 (e) Did you discuss whether the data you are using/curating contains personally identifiable
609 information or offensive content? [Yes] Please see Section 2. We have anonymized
610 the user ids in the submissions. The data samples are computer code for solving
611 context problems and in principle should not have any offensive content. While we
612 cannot guarantee that there is no personal information (e.g. name of a person) and
613 potential offensive content, but we have made every possible effort to minimize any
614 such possibility.

615 5. If you used crowdsourcing or conducted research with human subjects...

- 616 (a) Did you include the full text of instructions given to participants and screenshots, if
617 applicable? [N/A] We did not use crowdsourcing or conduct research with human
618 subjects
- 619 (b) Did you describe any potential participant risks, with links to Institutional Review
620 Board (IRB) approvals, if applicable? [N/A] We did not use crowdsourcing or conduct
621 research with human subjects
- 622 (c) Did you include the estimated hourly wage paid to participants and the total amount
623 spent on participant compensation? [N/A] We did not use crowdsourcing or conduct
624 research with human subjects

625 **A Additional details about CodeNet**

626 **A.1 URL**

627 [https://developer.ibm.com/technologies/artificial-intelligence/data/project-](https://developer.ibm.com/technologies/artificial-intelligence/data/project-codenet/)
628 codenet/ is the landing page of the dataset. It contains links to download the full dataset and the
629 benchmarks datasets (similar to POJ-104) in C++, Python and Java, which users can use to perform
630 similarity and classification experiments.

631 https://github.com/IBM/Project_CodeNet is the link to the Project CodeNet repository,
632 which contains software that supports and complements the CodeNet dataset. There are productivity
633 tools to aggregate codes samples based on user criteria and pre-processing tools to transform code
634 samples into sequence of tokens, simplified parse trees and code graphs. The repository also contains
635 notebooks that illustrate the usage of some of the tools and source code and scripts we used to perform
636 the experiments in the paper.

637 The URLs are all accessible to the general public.

638 **A.2 Author statement**

639 IBM represents and warrants it is the original author of the dataset and has the right to re-publish
640 associated third-party code under open source license terms. IBM further represents and warrants it
641 has the authority to grant the rights and licenses (CDLA Permissive v2.0) associated with the dataset
642 to third parties.

643 **A.3 Hosting and maintenance plan**

644 The CodeNet dataset is hosted under the IBM Data Asset eXchange (DAX) platform, which is an
645 online hub open to IBM and external developers and data scientists to find free and open data sets
646 under open data licenses. For developers, DAX offers a trusted source for open data sets for artificial
647 intelligence (AI). These data sets are ready to use in enterprise AI applications and are supplemented
648 with relevant notebooks and tutorials. DAX was launched in 2019 and maintained by the Center for
649 Open-Source Data & AI Technologies (CODAIT) team, who has been working on steadily adding
650 new data sets to the exchange, as well as resources that help explore these data sets.

651 The Project CodeNet repository is hosted under [github.com/IBM](https://github.com/IBM/Project_CodeNet) and is maintained by the Project
652 CodeNet team in IBM Research.

653 **A.4 How to read the CodeNet dataset**

654 The data and metadata are organized in a rigorous directory structure. The top level Project_CodeNet
655 directory contains several sub-directories: `data`, `metadata`, `problem_descriptions`, and
656 `derived`. The code samples or submissions reside under the `data` directory. The `data` directory
657 is organized as `(problem_id)/(language)/(submission)`, so the file path `data/p00023/C++/s006384060.cpp`
658 denotes a submission to problem p00023 in C++ with id s006384060. Detailed
659 statement of the problems can be found in `problem_descriptions/(problem_id).html`. The
660 meta data for the dataset is contained in the `metadata` directory. `metadata/problem_list.csv`
661 contains metadata for all the problems in the dataset, which is summarized in Table 7. `metadata/`
662 `(problem_id).csv` contains the metadata for all the submissions to problem `problem_id`, which is
663 described in Table 8. Each submission comes with cpu time, memory usage and status with possible
664 values described in Table 9. The `derived` directory contains information derived from the dataset,
665 such as near-duplicate information for submissions to specific languages, token sequences for code
666 samples, and information on identical problems.

667 **A.5 Long term preservation**

668 The dataset is hosted in the IBM Data Asset eXchange and is stored on IBM Cloud. Project CodeNet
669 is IBM's long term research effort to encourage open innovation at the intersection of AI and Software
670 Engineering. IBM has demonstrated a sustained commitment to open source innovation and the
671 CodeNet dataset and repository will be maintained and enhanced as long as is needed.

Table 7: Metadata at the dataset level

name of column	data type	unit	description
id	string	none	unique anonymized id of the problem
name	string	none	short name of the problem
dataset	string	none	original dataset, AIZU or AtCoder
time_limit	int	millisecond	maximum time allowed for a submission
memory_limit	int	KB	maximum memory allowed for a submission
rating	int	none	rating, i.e., difficulty of the problem
tags	string	none	list of tags separated by ";" not used
complexity	string	none	degree of difficulty of the problem; not used

Table 8: Metadata at the problem level

name of column	data type	unit	description
submission_id	string	none	unique anonymized id of the submission
problem_id	string	none	anonymized id of the problem
user_id	string	none	anonymized user id of the submission
date	int	seconds	date and time of submission in the Unix timestamp format (seconds since the epoch)
language	string	none	mapped language of the submission (ex: C++14 ->C++)
original_language	string	none	original language specification
filename_ext	string	none	extension of the filename that indicates the programminglanguage used
status	string	none	acceptance status, or error type
cpu_time	int	millisecond	execution time
memory	int	KB	memory used
code_size	int	bytes	size of the submission source code in bytes
accuracy	string	none	number of tests passed; *Only for AIZU

672 A.6 License

673 The dataset is distributed under the CDLA Permissive v2.0 license (<https://github.com/>
 674 Community-Data-License-Agreements/Working-Drafts/blob/main/CDLA-Permissive-
 675 2.0.md and <https://www.linuxfoundation.org/category/press-release/linux->
 676 foundation-press-release/page/2/) The repository is under the Apache License 2.0
 677 (<https://www.apache.org/licenses/LICENSE-2.0>).

678 A.7 Persistent dereferenceable identifier

679 The DOI for the Project CodeNet code repository is 10.5281/zenodo.4814770.

Table 9: All the possible status values

status	abbreviation	numeric code
Compile Error	CE	0
Wrong Answer	WA	1
Time Limit Exceeded	TLE	2
Memory Limit Exceeded	MLE	3
Accepted	AC	4
Judge Not Available	JNA	5
Output Limit Exceeded	OLE	6
Runtime Error	RE	7
WA: Presentation Error	PE	8
Waiting for Judging	WJ	
Waiting for Re-judging	WR	
Internal Error	IE	
Judge System Error		

680 B Datasheet

681 B.1 Motivation

682 1. For what purpose was the dataset created?

683 The CodeNet dataset provides a very large dataset of software source code written in a diversity of
 684 programming languages to drive algorithmic innovations in AI for code tasks like: code translation,
 685 code similarity, code classification, code search etc.

686 2. Who created this dataset (e.g. which team, research group) and on behalf of which entity 687 (e.g. company, institution, organization)?

688 The CodeNet dataset is created by a team of scientists at IBM Research and MIT-IBM Watson AI
 689 Lab comprising Ruchir Puri, David S. Kung, Geert Janssen, Wei Zhang, Giacomo Domeniconi,
 690 Vladimir Zolotov, Julian Dolby, Jie Chen, Mihir Choudhury, Lindsey Decker, Veronika Thost,
 691 Luca Buratti, Saurabh Pujar, Shyam Ramji, Ulrich Finkler, Susan Malaika, and Frederick Reiss.

692 3. What support was needed to make this dataset?

693 Project CodeNet is a research project within the IBM Research Division, so it is funded by the
 694 IBM Corporation.

695 B.2 Composition

696 1. What are the instances?(that is, examples; e.g., documents, images, people, countries) Are 697 there multiple types of instances? (e.g., movies, users, ratings; people, interactions between 698 them; nodes, edges)

699 The dataset consists of computer programs that are submissions to online judging sites and their
 700 accompanying metadata. CodeNet does not have multiple types of instances.

701 2. How many instances are there in total (of each type, if appropriate)?

702 The dataset comprises 13,916,868 submissions, divided into 4053 problems (of which 5 are empty).
 703 Of the submissions 53.6% (7,460,588) are accepted, 29.5% are marked as wrong answer and the
 704 remaining suffer from one of the possible rejection causes. The data contains submissions in 55
 705 different languages, although 95% of them are coded in the six most common languages (C++,
 706 Python, Java, C, Ruby, C#). C++ is the most common language with 8,008,527 submissions (57%
 707 of the total) of which 4,353,049 are accepted.

708 3. What data does each instance consist of ? “Raw” data (e.g., unprocessed text or images)? 709 Features/attributes?

710 The data are files of software programs as is. The character encoding of each file is UTF-8.

711 4. Is there a label or target associated with each instance? If so, please provide a description.

712 Yes, each instance of data (file) has associated metadata that may be interpreted as labels. The
 713 problem that a certain instance (code sample) intends to solve may be used as a label in classi-
 714 fication and similarity. The acceptance status of each code sample, the CPU time, and memory
 715 footprint can also be used as labels.

716 5. **Is any information missing from individual instances? If so, please provide a description,**
717 **explaining why this information is missing (e.g., because it was unavailable). This does not**
718 **include intentionally removed information, but might include, e.g., redacted text.**

719 Some metadata values might not be available for all instances. This can be attributed to the source
720 not (or incorrectly) providing the metadata.

721 6. **Are relationships between individual instances made explicit (e.g., users' movie ratings,**
722 **social network links)? If so, please describe how these relationships are made explicit.**

723 Relationships between instances are explicitly available in the provided metadata. As an example,
724 multiple instances (code submissions) by the same person can be found by scanning the metadata
725 for that person's (anonymized) id number. All submissions to a particular problem id are to be
726 found in a single metadata CSV file.

727 7. **Are there recommended data splits (e.g., training, development/validation, testing)? If so,**
728 **please provide a description of these splits, explaining the rationale behind them.**

729 Data splits are left to the discretion of the user, since CodeNet can be used for a wide variety of
730 use cases. No such division is made in the dataset.

731 8. **Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide**
732 **a description.**

733 It all depends on how errors, noise and redundancies are defined. There are probably minor errors
734 in the metadata directly attributable to the source: some non-Accepted programs are identified
735 with the wrong language, probably caused by a programmer making a wrong selection while
736 submitting his or her work. Some run-time data for incorrect programs are listed as a negative
737 number. It happens that some user or users submit the same program (data instance) multiple
738 times to the same or different problem tasks.

739 9. **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**
740 **websites, tweets, other datasets)? If it links to or relies on external resources, a) are there**
741 **guarantees that they will exist, and remain constant, over time; b) are there official archival**
742 **versions of the complete dataset (i.e., including the external resources as they existed at**
743 **the time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associ-**
744 **ated with any of the external resources that might apply to a future user? Please provide**
745 **descriptions of all external resources and any restrictions associated with them, as well as**
746 **links or other access points, as appropriate.**

747 The dataset is self-contained.

748 10. **Does the dataset contain data that might be considered confidential (e.g. data that is pro-**
749 **tected by legal privilege or by doctor-patient confidentiality, data that includes the content**
750 **of individuals' non-public communications)? If so, please provide a description.**

751 No. Any personal information that is available in the metadata at the source websites is anonymized
752 in the dataset. However, it is possible that names or handles of persons still being present in the
753 source code instances themselves as variable, class or function names.

754 11. **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threat-**
755 **ening, or might otherwise cause anxiety? If so, please describe why.**

756 We have done some filtering. In one case, the programming language name is offensive, so we
757 renamed it. It might be possible that people used offensive language in naming a variable in the
758 program, but we have made every possible effort to minimize any such possibility.

759 12. **Does the dataset relate to people?**

760 No.

761 **B.3 Collection Process**

762 1. **How was the data associated with each instance acquired? Was the data directly observable**
763 **(e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly**
764 **inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or**
765 **language)? If data was reported by subjects or indirectly inferred/derived from other data,**
766 **was the data validated/verified? If so, please describe how.**

767 The data are acquired from publicly accessible on-line judging websites. We used the AIZU
768 (<https://judge.u-aizu.ac.jp/>) and AtCoder (<https://atcoder.jp/>) online judging sites. The data are
769 accessible and observable by clicking specific url links.

- 770 2. **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or**
771 **sensor, manual human curation, software program, software API)? How were these mech-**
772 **anisms or procedures validated?**
- 773 Some of the data was available as archived zip files or a REST API for download, otherwise a
774 webpage scraper tool was used to retrieve the data (while observing any throttles on bandwidth).
775 No verification beyond mere manual inspection was applied to the downloaded data.
- 776 3. **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., determin-**
777 **istic, probabilistic with specific sampling probabilities)?**
- 778 No specific strategy: as much data as was available at the time (2020) was downloaded.
- 779 4. **Who was involved in the data collection process (e.g., students, crowdworkers, contractors)**
780 **and how were they compensated (e.g., how much were crowdworkers paid)?**
- 781 There were no third-party participants in the data collection.
- 782 5. **Over what timeframe was the data collected? Does this timeframe match the creation time-**
783 **frame of the data associated with the instances (e.g., recent crawl of old news articles)?**
784 **If not, please describe the timeframe in which the data associated with the instances was**
785 **created.**
- 786 The data was collected in 2020 and the code samples might go back to a decade ago. The dataset
787 was first published on May 5, 2021.
- 788 6. **Were any ethical review processes conducted (e.g. by an institutional review board)? If so,**
789 **please provide a description of these review processes, including the outcomes, as well as a**
790 **link or other access point to any supporting documentation.**
- 791 No. The dataset was examined by IBM Corporation lawyers for suitability of public disclosure.
- 792 7. **Does the dataset relate to people? If not, you may skip the remainder of the questions in**
793 **this section.**
- 794 Only as far as the fact that the data instances are created/written by people.
- 795 8. **Did you collect the data from the individuals in question directly, or obtain it via third**
796 **parties or other sources (e.g. websites)?**
- 797 The data was collected indirectly from submitters via online judging websites.
- 798 9. **Did the individuals in question consent to the collection and use of their data? If so, please**
799 **describe (or show with screenshots or other information) how consent was requested and**
800 **provided, and provide a link or other access point to, or otherwise reproduce, the exact**
801 **language to which the individuals consented.**
- 802 They did consent directly to the respective online judging sites that we used as source. See e.g.
803 https://onlinejudge.u-aizu.ac.jp/term_of_use.

804 **B.4 Data Preprocessing/Cleaning**

- 805 1. **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,**
806 **tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, pro-**
807 **cessing of missing values)? If so, please provide a description. If not, you may skip the**
808 **remainder of the questions in this section.**
- 809 Minor processing of the data instances was performed mainly to make all file name extensions
810 uniform and make sure the character encoding is UTF-8, all line endings adhere to the UNIX
811 standard (a single linefeed character), and any byte-order marks (BOM) are removed.
- 812 All metadata was carefully examined, anonymized where necessary and corrected or updated
813 when possible (e.g. the file size in bytes is part of the metadata and needed updating).
- 814 2. **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**
815 **support unanticipated future uses)? If so, please provide a link or other access point to the**
816 **“raw” data.**
- 817 The raw data is saved but considered not to be part of the published Project CodeNet dataset.
- 818 3. **Is the software used to preprocess/clean/label the instances available? If so, please provide**
819 **a link or other access point.**
- 820 Some of the software (mostly bash scripts) are available in our github https://github.com/IBM/Project_CodeNet.

822 4. **Does this dataset collection/processing procedure achieve the motivation for creating the
823 dataset stated in the first section of this datasheet? If not, what are the limitations?**

824 Yes. This dataset and its derived benchmark datasets offer the scale, diversity, and quality to drive
825 research in applying AI techniques to code.

826 **B.5 Uses**

827 1. **Has the dataset been used for any tasks already? If so, please provide a description.**

828 Yes. Several baseline experiments on code classification and similarity have been performed and
829 documented in the paper.

830 2. **Is there a repository that links to any or all papers or systems that use the dataset? If so,
831 please provide a link or other access point.**

832 Yes. https://github.com/IBM/Project_CodeNet.

833 3. **What (other) tasks could the dataset be used for?**

834 The rich metadata and diversity open Project CodeNet to a plethora of uses cases. The problem-
835 submission relationship in Project CodeNet corresponds to type-4 similarity and can be used for
836 code search and clone detection. The code samples in Project CodeNet are labeled with their
837 acceptance status and we can explore AI techniques to distinguish correct codes from problematic
838 ones. Project CodeNet's metadata also enables the tracking of how a submission evolves from
839 problematic to accepted, which could be used for exploring automatic code correction. A large
840 number of code samples come with inputs so that we can execute the codes to extract the CPU run
841 time and memory footprint, which can be used for regression studies and predictions. Given its
842 wealth of programs written in a multitude of languages, Project CodeNet may serve as a valuable
843 benchmark dataset for source-to-source translation.

844 4. **Is there anything about the composition of the dataset or the way it was collected and pre-
845 processed/cleaned/labeled that might impact future uses? For example, is there anything
846 that a future user might need to know to avoid uses that could result in unfair treatment
847 of individuals or groups (e.g. stereotyping, quality of service issues) or other undesirable
848 harms (e.g. financial harms, legal risks)? If so, please provide a description. Is there anything
849 a future user could do to mitigate these undesirable harms?**

850 No.

851 5. **Are there tasks for which the dataset should not be used? If so, please provide a description.**

852 No.

853 **B.6 Dataset Distribution**

854 1. **Will the dataset be distributed to third parties outside of the entity (e.g. company, insti-
855 tution, organization) on behalf of which the dataset was created? If so, please provide a
856 description.**

857 Yes, the dataset will be distributed to the general public.

858 2. **When will the dataset be released/first distributed? What license (if any) is it distributed
859 under?** The dataset was released in May 2021 under the the CDLA Permissive v2.0 li-
860 cence [https://github.com/Community-Data-License-Agreements/Working-Drafts/
861 blob/main/CDLA-Permissive-2.0.md](https://github.com/Community-Data-License-Agreements/Working-Drafts/blob/main/CDLA-Permissive-2.0.md).

862 3. **How will the dataset be distributed (e.g. tarball on website, API, GitHub)? Does the dataset
863 have a digital object identifier (DOI)?**

864 The dataset is made available as a downloadable gzipped tar file here: [https://developer.
865 ibm.com/technologies/artificial-intelligence/data/project-codenet/](https://developer.ibm.com/technologies/artificial-intelligence/data/project-codenet/). There is
866 no DOI yet.

867 4. **Will the dataset be distributed under a copyright or other intellectual property (IP) license,
868 and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU,
869 and provide a link or other access point to, or otherwise reproduce, any relevant licensing
870 terms or ToU, as well as any fees associated with these restrictions.**

871 The dataset is made available under the CDLA Permissive v2.0 licence [https://github.com/Community-Data-License-Agreements/Working-Drafts/
873 blob/main/CDLA-Permissive-2.0.md](https://github.com/Community-Data-License-Agreements/Working-Drafts/blob/main/CDLA-
872 Permissive-2.0.md).

874 **5. Have any third parties imposed IP-based or other restrictions on the data associated with**
875 **the instances? If so, please describe these restrictions, and provide a link or other access**
876 **point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated**
877 **with these restrictions.**

878 No, not as far as we know.

879 **6. Do any export controls or other regulatory restrictions apply to the dataset or to individual**
880 **instances? If so, please describe these restrictions, and provide a link or other access point**
881 **to, or otherwise reproduce, any supporting documentation.**

882 No. These code samples are solutions to pedagogical programming problems at the high school
883 and beginning college level and should not be subject to export control.

884 **B.7 Dataset Maintenance**

885 **1. Who is supporting/hosting/maintaining the dataset?**

886 International Business Machines corporation.

887 **2. How can the owner/curator/manager of the dataset be contacted (e.g. email address)?**

888 The users can create an issue on our github or contact any of the listed authors.

889 **3. Is there an erratum? If so, please provide a link or other access point.**

890 No.

891 **4. Will the dataset be updated (e.g. to correct labeling errors, add new instances, delete in-**
892 **stances)? If so, please describe how often, by whom, and how updates will be communicated**
893 **to users (e.g. mailing list, GitHub)?**

894 Yes, there are plans to add new instances to the dataset, in the next six months to a year's time
895 frame. The update will be performed by IBM and communicated through the github.

896 **5. If others want to extend/augment/build on this dataset, is there a mechanism for them to**
897 **do so? If so, please provide a description. Will these contributions be validated/verified? If**
898 **so, please describe how. If not, why not? Is there a process for communicating/distributing**
899 **these contributions to other users? If so, please provide a description.**

900 There is no such mechanism yet, but it is under consideration. Interested parties are invited to
901 consider contacting the authors or creating an issue to that effect in our github.

902 **C Further information of CodeNet**

903 Table 10 summarizes the metadata available for each code submission to a problem. Figure 4 gives
 904 the distributions of problems based on number of submissions received.

Table 10: Submission metadata.

column	unit/example	description
submission_id	s[0-9]{9}	anonymized id of submission
problem_id	p[0-9]{5}	anonymized id of problem
user_id	u[0-9]{9}	anonymized user id
date	seconds	date and time of submission
language	C++	consolidated programming language
original_language	C++14	original language
filename_ext	.cpp	filename extension
status	Accepted	acceptance status, or error type
cpu_time	millisecond	execution time
memory	kilobytes	memory used
code_size	bytes	source file size
accuracy	4/4	passed tests (AIZU only)

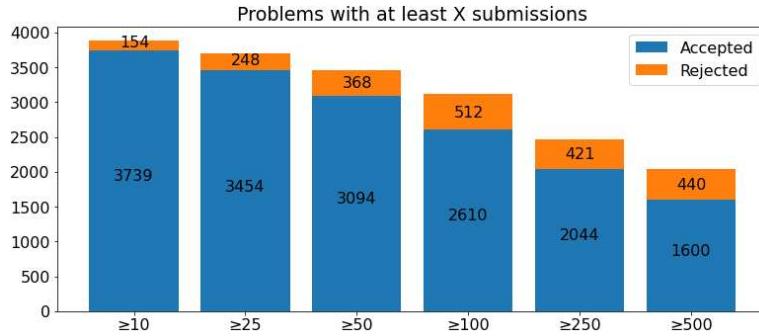


Figure 4: Number of problems providing at least X submissions. The bars show both the numbers of accepted submissions (blue) and rejected submissions (orange).

905 **D Details of Experiments on Code Classification**

906 **D.1 MLP with Bag of Tokens**

907 One of the simplest representations of a code sample is a bag of tokens. Here, the code sample is
 908 represented by a vector of relative frequencies of token occurrences in the source code. The vector is
 909 computed by the following steps:

- 910 1. Convert a given source code into a sequence of tokens using a tokenizer (i.e., lexical analyzer).
 911 2. From this sequence, remove the tokens considered not useful for code classification.
 912 3. Count the number of each token type in the reduced sequence and form a vector of counts.
 913 4. Normalize the vector with respect to L2 norm.

914 We do not use all tokens available in the grammar of the programming language. Only some operators
 915 and keywords are used. All identifiers, comments and literals are ignored. We also ignore some
 916 operators and many keywords that in our opinion provide no significant information on the algorithm
 917 the source code implements.

918 The vector representing a bag of tokens has the same length for every code sample, which makes
 919 it convenient for processing with a neural network. The vector is usually short, which makes
 920 training of a neural network fast. However, in a bag-of-tokens representation, information about the
 921 number of occurrences and position of each token is lost. Hence, the accuracy of a classifier using a
 922 bag-of-tokens representation is rather limited.

923 Table 11 provides results of code classification of all four benchmarks. The columns give the
 924 benchmark name, the test accuracy, the number of training epochs, the run time of each epoch, and
 925 the number of token types considered. All networks are implemented using Keras API of TensorFlow
 926 machine learning tool. Training is performed on a single V100 GPU, using Adam optimizer with
 927 learning rate 1e-3, and batches of 32 samples. In each experiment, 20% of the samples are used for
 928 testing, while the rest are split in 4:1 for training and validation, respectively.

Table 11: Code classification by MLP with bag of tokens.

Benchmark dataset	Accuracy %%	Number epochs	Run time sec/epoch	Number tokens
Java250	71.00±0.29	30	2	81
Python800	67.80±0.15	22	7	71
C++1000	68.26±0.21	20	14	56
C++1400	64.50±0.13	17	12	56

929 Figure 5 shows the neural network used for solving the classification problem for the C++1400
 930 benchmark. The neural networks used for classification of other benchmarks are similar to this one.
 931 As we see in Table 11 their performance is quite similar.

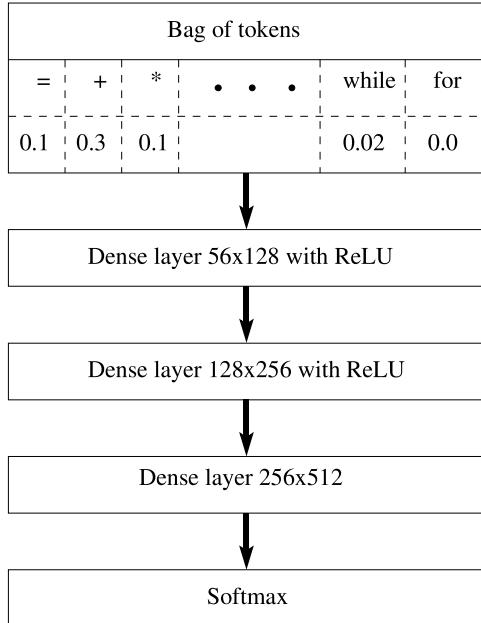


Figure 5: MLP architecture for code classification.

932 From Table 11 we see that training is rather fast, the reason being that the network is simple. In
 933 spite of simplicity, this neural network performs very well. The 64.50±0.13% test accuracy for
 934 C++1400 benchmark dataset is significantly better than the potential 0.071% accuracy of random
 935 guess. It indicates that the relative frequencies of source code tokens provide sufficient information
 936 for classifying code.

937 D.2 CNN with Token Sequence

938 The sequence-of-tokens representation retains more information of a code sample than the bag-of-
 939 tokens representation. For our experiments on code classification, we use the same set of tokens that
 940 is used in the above bag-of-tokens approach. Similarly, we omit all comments and identifiers.

941 Table 12 shows results of code classification on all four benchmarks by using the sequence-of-tokens
 942 representation. The columns give the benchmark name, the test accuracy, the number of training
 943 epochs, the run time of each epoch, and the number of token types considered. All networks are
 944 implemented using Keras API of TensorFlow machine learning tool. The training is performed on

Table 12: Code classification by CNN with token sequence.

Benchmark dataset	Accuracy %	Number epochs	Run time sec/epoch	Number tokens
Java250	89.52±0.59	810	10	81
Python800	87.46±0.25	504	26	71
C++1000	93.96±0.18	235	59	56
C++1400	93.71±0.18	334	60	56

945 four V100 GPUs, using Adam optimizer in data parallel mode with learning rate 1e-3, and batches of
946 512 samples. In each experiment, 20% of the samples are used for testing, while the rest are split in
947 4:1 for training and validation, respectively.

948 We have experimented with several types of neural networks. Figure 6 shows the neural network
949 we choose for the C++1400 benchmark. It is a multi-layer convolutional neural network. It uses
950 categorical encoding of source code tokens. For batching, the sequences of tokens are padded with
951 zeros.

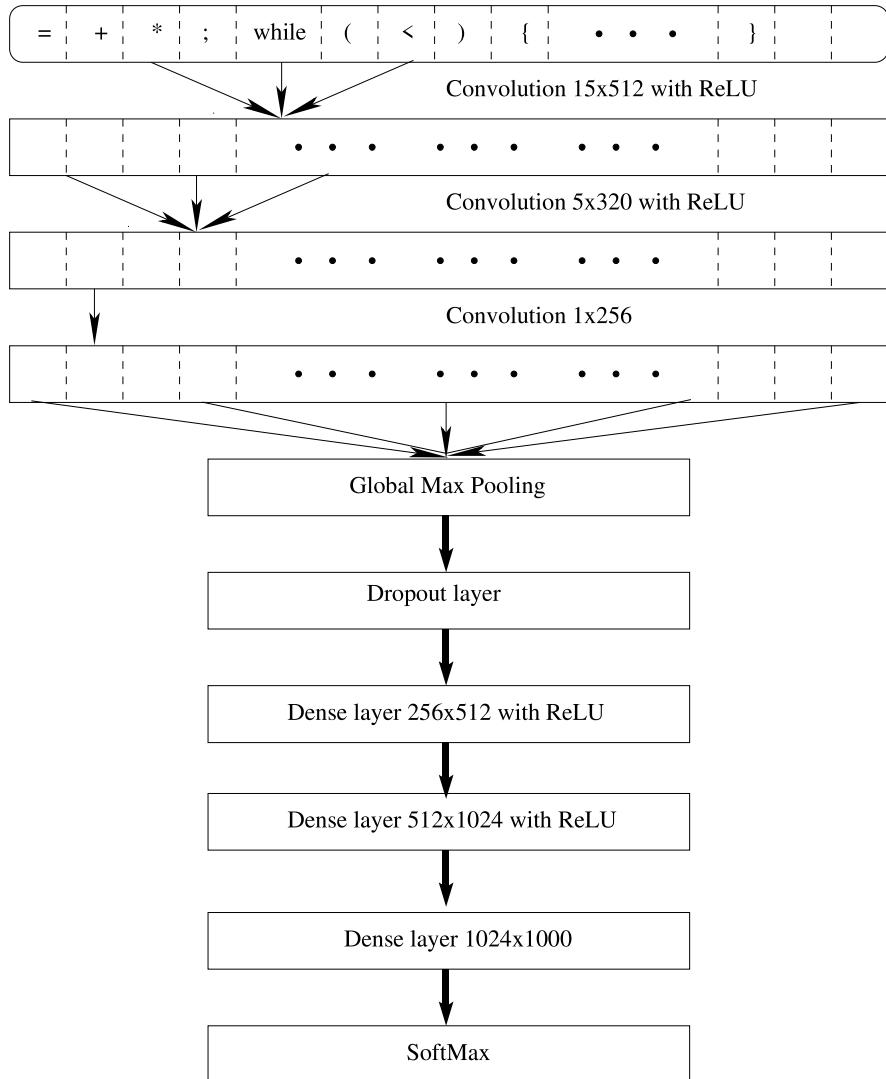


Figure 6: CNN architecture for code classification.

952 Using this network we get a test accuracy 93.71±0.18% for C++1400 benchmark dataset, which is
953 significantly better than the accuracy shown by the bag-of-tokens approach. The neural networks

954 used for classification of other benchmarks are similar to the one shown in Figure 6. As we see in
955 Table 12, their performance is similar.

956 D.3 C-BERT with Token Sequence

957 The sequence-of-tokens representation can be used with other neural networks of increasing capacity.
958 We build a C-BERT model (a transformer model introduced in [33]) by pre-training on 10,000 top
959 starred GitHub open source projects written in C, where we use Clang C tokenizer and Sentencepiece
960 to tokenize the pre-training data. The C-BERT model is then fine tuned on each classification
961 benchmark. Additionally, we experiment with the POJ-104 dataset, which contains code examples in
962 C and C++.

963 C-BERT achieves appealing results on binary classification and vulnerability detection with C source
964 code [10, 49]. However, it has not been used on multiclass classification tasks or with other languages
965 such as C++, Java, and Python. Because we use sub-word tokenization and different programming
966 languages share common tokens, we could apply the C-BERT model directly on the benchmarks.

967 After pretraining, we fine tune the model for five epochs on each benchmark, with a batch size 32 and
968 learning rate 2e-5. The fine-tuning was done on two V100 GPUs and it took 30 minutes to four hours,
969 depending on the size of the dataset. The sub-word vocabulary size is 5,000. Contexts larger than
970 512 tokens were truncated.

971 Table 13 summarizes the accuracies C-BERT achieves on the four CodeNet benchmarks as well as the
972 POJ-104 dataset. C-BERT achieves high accuracy and performs the best on Java and Python.

Table 13: C-BERT results (accuracy, in %) for code classification.

	POJ-104	C++1000	C++1400	Java250	Python800
C-BERT	98.41±0.01	93.79±0.01	91.83±0.06	97.40±0.19	97.09±0.18

973 The relatively low performance on C++ benchmarks is possibly related to the idiosyncrasies of the
974 dataset and certain programming practices. Manual inspection suggests that lack of detailed variable
975 names in C++ hurts the performance of the model, in problems appearing similar and having similar
976 solutions. Removing one of the similar problems improves the model performance on the other
977 problem. Moreover, one programming practice which could potentially confuse the models is that
978 certain C++ users copied common constants (e.g., pi and epsilon) and data structures (e.g., enums) to
979 all solutions they submitted. In many cases, these duplicate contents were not even used. We did not
980 observe such practices in Python and Java.

981 D.4 GNN with SPT

982 We experiment with four types of GNNs with SPT-based graph representations of the source code:
983 the Graph Convolutional Network (GCN) [34], the Graph Isomorphism Network (GIN) [35], and a
984 virtual-node-included variant for each (denoted by -V). The variant adds a virtual node to the graph
985 to enhance graph message passing [36]. We use the Adam optimizer with learning rate 1e-3 for
986 training. All GNN models have five layers. We have experimented with more than 5 layers (i.e., 8
987 and 10), however deeper GNNs do not improve performance, as deeper GNNs might suffer from
988 the over-smoothing problem (i.e., node features become less distinguishable after many rounds of
989 message passing) [50].

990 We conduct 6/2/2 random split for each of the 4 benchmarks: i.e., 60% training data, 20% testing
991 data, and 20% validation data. We run five folds for each benchmark with early stop "patience"
992 set 20 (i.e., stop only when validation loss has not decreased in the past 20 epochs). Our model
993 training typically converges within 200 epochs in a 1-fold run. We modified OGB [51] code-base with
994 PyTorch Geometric [52] back-end over PyTorch 1.6.0 [53] to run our experiments. The experiments
995 are conducted on one NVIDIA V100 GPU. For large benchmarks such as C++1000 and C++1400, it
996 takes about 1 week to finish a 5-fold run. We summarize model accuracy, training time over 5-folds,
997 and training epochs over 5-folds in Table 14. As we can see, adding a virtual node improves GNN
998 performance (both GCN and GIN). Overall, GIN and its variants work better than GCN and its
999 variants, likely due to the fact that GIN theoretically generalizes the Weisfeiler-Lehman Isomorphism
1000 Test and achieves maximum expressive power among GNNs [54].

1001 For the detailed model, hyper-parameter setup, data splits and etc, please refer to https://github.com/IBM/Project_CodeNet/tree/main/model-experiments/gnn-based-experiments.
 1002

Table 14: GNN (SPT) results for code classification. Each task trains over 5-folds with early stopping patience parameter set as 20. We record test accuracy (with standard deviation), total training time over 5 folds, and total training epochs over 5 folds.

	Java250	Python800	C++1000	C++1400
GCN	92.70±0.25 10.55 hrs 411 epochs	93.82±0.16 14.50 hrs 219 epochs	95.76±0.12 47.96 hrs 228 epochs	95.26±0.13 67.34 hrs 310 epochs
	93.02±0.81 12.50 hrs 419 epochs	94.30 ±0.15 23.02 hrs 325 epochs	96.09±0.17 61.55 hrs 287 epochs	95.73±0.07 71.85 hrs 358 epochs
	93.26±0.23 19.80 hrs 513 epochs	94.17±0.19 41.67 hrs 496 epochs	96.34±0.15 116.67 hrs 441 epochs	95.95±0.13 133.50 hrs 502 epochs
GIN	92.77±0.66 26.25 hrs 656 epochs	94.54±0.12 51.67 hrs 570 epochs	96.64±0.10 142.25 hrs 496 epochs	96.36±0.10 208.47 hrs 678 epochs
	92.77±0.66 26.25 hrs 656 epochs	94.54±0.12 51.67 hrs 570 epochs	96.64±0.10 142.25 hrs 496 epochs	96.36±0.10 208.47 hrs 678 epochs
	92.77±0.66 26.25 hrs 656 epochs	94.54±0.12 51.67 hrs 570 epochs	96.64±0.10 142.25 hrs 496 epochs	96.36±0.10 208.47 hrs 678 epochs

1003 E Details of Experiments on Code Similarity

1004 E.1 MLP with Bag of Tokens

1005 For experiments on code similarity analysis, we use the same bag of tokens as for code classification.
 1006 The input to the neural network is constructed by concatenating two bags of tokens, one for each
 1007 source code file.

1008 Table 15 provides results of code similarity analysis on all four benchmarks. The columns give the
 1009 benchmark name, the test accuracy, the number of training epochs, the number of samples in each
 1010 epoch, the run time of each epoch, the number of token types considered, and the number of test
 1011 samples. All networks are implemented using Keras API of TensorFlow machine learning tool. The
 1012 training is performed on a single V100 GPU, using Adam optimizer with learning rate 1e-3, and
 1013 batches of 256 samples.

Table 15: Similarity analysis by MLP with bag of tokens.

Benchmark dataset	Accuracy %%	Number epochs	Size of epoch	Run time sec/epoch	Number tokens	N test samples
Java250	81.80±0.06	20	4,096,000	21	81	512,000
Python800	86.61±0.08	94	4,096,000	24	71	512,000
C++1000	85.82±0.05	64	4,096,000	21	56	512,000
C++1400	86.54±0.07	64	4,096,000	22	56	512,000

1014 Figure 7 shows the neural network used for code similarity analysis on the C++1400 benchmark. The
 1015 neural networks used for code similarity analysis on other benchmarks are similar to this one. As we
 1016 see in Table 15, their accuracy is similar.

1017 As we see in Table 15, the model accuracy is rather modest (<87%) for all benchmark datasets, which
 1018 is not very high for a binary classification problem of a fully balanced dataset. Obviously, the bag of
 1019 tokens is too primitive and misses many important details necessary for identifying similarity.

1020 E.2 Siamese Network with Token Sequence

1021 For experiments on code similarity, we use the same sequence of tokens as for code classification.
 1022 The neural network has two inputs, one for each source code file. After experimenting with various
 1023 neural network architectures, we select the siamese network for its good performance.

1024 Table 16 provides results of code similarity analysis on all four benchmarks. The columns give the
 1025 benchmark name, the test accuracy, the number of training epochs, the number of samples in each

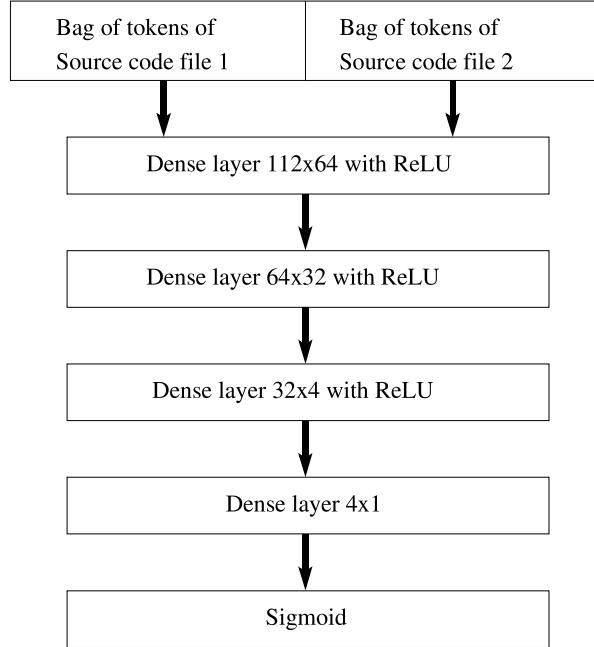


Figure 7: MLP architecture for similarity analysis.

1026 epoch, the run time of each epoch, the number of token types considered, and the number of test
 1027 samples. All networks are implemented using Keras API of TensorFlow machine learning tool. The
 1028 training is performed on four V100 GPUs, using Adam optimizer in data parallel mode with learning
 1029 rate $1e-3$, and batches of 512 samples.

Table 16: Similarity analysis by Siamese network with token sequence.

Benchmark dataset	Accuracy %	Number epochs	Size of epoch	Run time sec/epoch	Number tokens	N test samples
Java250	89.70 ± 0.18	29	51,200	114	75	512,000
Python800	94.67 ± 0.12	110	64,000	89	71	512,000
C++1000	96.19 ± 0.08	123	64,000	89	56	512,000
C++1400	96.56 ± 0.07	144	64,000	96	56	512,000

1030 The neural network for the C++1400 benchmark is depicted in Figure 8. The siamese parts of the
 1031 network have the same structure and share all their weights. If the inputs are identical, so are the
 1032 outputs. Therefore, by construction, the network guarantees detecting similarity of identical source
 1033 code samples. The outputs of the siamese parts are compared by computing the absolute difference.

1034 The network shows $96.56 \pm 0.07\%$ test accuracy for C++1400 benchmark dataset. We consider this a
 1035 good result, especially considering that the token sequence ignores all identifiers, comments, and
 1036 many keywords. The neural networks used for code similarity analysis of other benchmarks are
 1037 similar to the one shown in Figure 8. As we see in Table 16, their accuracy is quite similar.

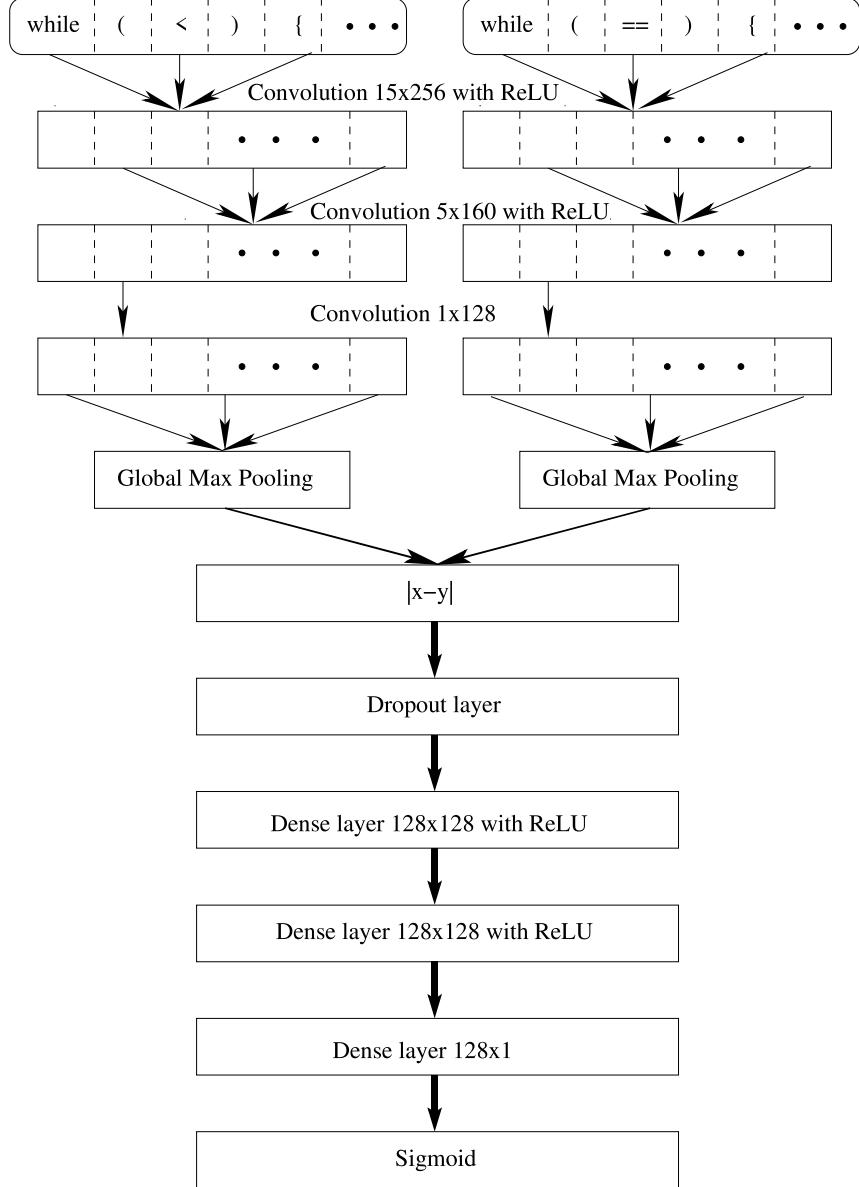


Figure 8: Siamese architecture for similarity analysis.

1038 **E.3 SPT-based experiments**

1039 Following MISIM [21], the train, validation, and test datasets for the SPT-based experiments draw
 1040 from entirely different problems. In our experiments, we use 50% problems for training, 25%
 1041 for validation, and 25% for test. The train, validation, and test split used for the experiments can
 1042 be found at [55]. Similarity scores in Table 5 and Table 6 report mean and standard deviation of
 1043 MAP@R [40] values evaluated with models trained using five random seeds. The models are trained
 1044 on a Xeon(R) CPU E5-2680 v4, 2.4GHz, 256 GiB memory using a NVIDIA V100 GPU. The SPTs
 1045 used in these experiments have nodes annotated with attributes derived by combining SPT features
 1046 (refer to Section 6), following the context-aware semantic structure (CASS) proposed in [21].

1047 AROMA experiments are performed using the implementation in MISIM's supplementary mate-
 1048 rial [23] and the input (SPTs) used for these experiments can be found at [55]. Due to the high
 1049 memory requirement for computing MAP@R on the test set of CodeNet benchmarks, we had to
 1050 reduce the feature set of AROMA. We estimate that AROMA results can improve by 10–25% when

1051 all features are used. AROMA is rule-based and no training is involved, hence we don't report mean
 1052 and standard deviation in Table 5. For each of the four datasets – Java250, Python800, C++1000,
 1053 C++1400 – MISIM's GNN model is trained for a total of 1000 epochs at a learning rate of 0.001
 1054 with Adam optimizer. Each epoch consists of 1000 iterations, and in each iteration, 16 problems
 1055 and 5 solutions per problem are randomly sampled, and all solution pairs are used for training as in
 1056 [21]. MISIM results for the four languages can be reproduced by downloading the MISIM code and
 1057 scripts [23] and using the provided CASS files [55] as input.
 1058 For the GMN experiments (row 2 and row 3 in Table 6), we adapt the implementation in [39] of
 1059 the GMN model [38] using SPTs [55] as graphs. We follow the recommendations in [38] for the
 1060 model configuration, as they produce the best and stable results in our experiments. Specifically,
 1061 we use 5 layers of propagation with weight sharing across layers, dot-product similarity for the
 1062 cross-graph attention mechanism, and GRU layer to update node embeddings from the propagation
 1063 scheme. For GMN training, given the large set of SPT pairs, we adopt an approach similar to [21] of
 1064 randomly sampling 16 problems with 5 solutions each. We use triplet loss with approximate hamming
 1065 similarity [38] for each sample, which is formed using a similar pair combined with a dissimilar SPT.
 1066 After every 100 iterations with a batch size of 64, another set of 16 problems and 5 solutions are
 1067 sampled randomly for a total of 150,000 iterations (1500 sampled sets). GMN results could improve
 1068 further with more training iterations. We use Adam optimizer with a learning rate of 1e-4 for training.
 1069 The first two rows of Table 6 compare similarity models trained on SPT graph structure only. The first
 1070 row in the table adapts the MISIM GNN model by masking the node labels to allow the model to learn
 1071 structural features only. The second row uses the GMN [38] model with cross-graph attention-based
 1072 matching for structural similarity using a node vector dimension of 32 and graph representation
 1073 dimension of 128.
 1074 For the GMN+MISIM node attributes experiment, row 3 in Table 6, we allow the GMN model to
 1075 learn features based on both node attributes and the SPT structure. Accordingly, we replace the node
 1076 encoder in the GMN, an MLP, with an embedding layer, for generating node feature vectors. We
 1077 explore different node feature vector dimensions, such as 64, 100, 128, and found 100 to produce
 1078 good results for the given number of training iterations. All other parameter settings remain the same
 1079 as the structure only GMN experiments from row 2 of Table 6. The GMN results can be reproduced
 1080 using the Java250 CASS files available at [55].
 1081 MAP@R score [40] is computationally expensive for GMN models because an embedding has to be
 1082 computed for all SPT pairs in the test set, and hence Table 6 reports results on smaller sampled test
 1083 sets.

1084 F Details of MLM Experiment

1085 Here we show how a masked language model (MLM) can be trained with CodeNet. We closely
 1086 follow the approach by Ankur Singh, documented in the blog [56]. The goal of the model is to infer
 1087 the correct token for an arbitrary masked-out location in the source text. We assume that in every text,
 1088 precisely one token is randomly masked. The original token at such position is then the golden label.
 1089 From each of the 1000 C++1000 problems, we randomly select 100 samples for training and another
 1090 100 for testing. Each C++ source file is tokenized into a vocabulary of 442 distinct tokens as
 1091 categorized in Table 17. For example, `while` is a keyword and `strlen` is a function literal.

Table 17: Token categories used for MLM.

Type	Count	Description
the keyword	95	all C++20 reserved words
the function	280	function names in common header files
the identifier	42	standard identifiers, like <code>stderr</code> , etc.
the punctuator	16	small set of punctuation symbols
# or ##	2	the C pre-processor symbols
0, 1	2	special case for these frequent constants
the token class	5	identifier, number, operator, character, string

1092 This code snippet:
 1093 `for (i = 0; i < strlen(s); i++) {}`

1094 will be tokenized to:
1095 `for (id = 0 ; id < strlen (id) ; id operator) { }`
1096 The tokenized source files are read into a pandas dataframe and processed by the Keras Text Vector-
1097 ization layer, to extract a vocabulary and encode all token lines into vocabulary indices, including the
1098 special “[mask]” token. Each sample has a fixed token length of 256. The average number of tokens
1099 per sample across the training set is 474. Short samples are padded with 0 and those that are too large
1100 are simply truncated.
1101 The model is trained with 100,000 samples in batches of 32 over five epochs, with a learning rate
1102 of 0.001 using the Adam optimizer. We evaluate the trained model on a test set of 100,000 samples.
1103 Each sample is pre-processed in the same way as the training samples and one token (never a padding)
1104 is arbitrarily replaced by the “[mask]” symbol. Then, a prediction is generated and the top 1 and top
1105 5 results are compared with the expected value. The achieved accuracies are top-1: 0.9104 (stddev:
1106 0.002) and top-5: 0.9935 (stddev: 0.0005).