

CODE-PROMPT: Evaluating Code-LLMs for their NLP Abilities via Code Aligned Prompts

Anonymous ACL submission

Abstract

Large language models pre-trained on code (Code-LLMs) have achieved remarkable performance on coding tasks. Despite their success of processing programming language, their application to natural language tasks has largely been limited to narrow tasks and ad-hoc model selection. In this paper, we present CODE-PROMPT, a comprehensive prompt-based framework to benchmark Code-LLMs over ten NLP tasks. We evaluate CODE-PROMPT on 13 Code-LLMs using code-aligned prompts and compare them with their parallel natural language LLMs (NL-LLMs) using natural language prompts.

Our results show that Code-LLMs perform on par with NL-LLMs across a range of NLP tasks, while producing more format-consistent generation with less redundancy. Our cross-lingual experiments further indicate that knowledge in Code-LLMs is well shared among different programming languages, whereas performance across different human languages still varies significantly. All our code and data will be made publicly available at [anonymous_url](#).

1 Introduction

The success of large language models (LLMs) (Touvron et al., 2023; OpenAI et al., 2024) across NLP tasks has demonstrated the power of large-scale pre-training on natural language. As prompting has become increasingly important (Liu et al., 2023; Long et al., 2025), model inputs (*prompts*) often include a task description, few-shot demonstrations, and a final example for inference (Brown et al., 2020). Following the trend of large-scale training, large models have been also pre-trained on massive code data for better code generation tasks (Wang et al., 2021; Rozière et al., 2023; Zhao et al., 2024). These Code-LLMs now power modern code assistants, such as Github Copilot¹ and Cursor². Besides ex-

ecutable code, Code-LLMs can also generate *natural language* elements in code repositories such as comment (Geng et al., 2024), docstring (Poudel et al., 2024), repository-level documentation (Luo et al., 2024) and Git commit messages (Zhang et al., 2024). In addition, Yang et al. (2024) and Wang et al. (2025) show that Code-LLMs can generate test cases for **algorithmic functions**, such as those found in programming challenges on LeetCode³.

With the rise of NLP-related projects in industry (Supriyono et al., 2024), engineers and data scientists increasingly need to prepare test examples with expected outputs for **NLP functions**. In software engineering, such test cases are essential for unit tests and often appear in (Python) docstrings; in LLM-driven projects, developers design prompts and require such test cases inserted as few-shot demonstrations to guide model behaviour. A natural solution is to leverage Code-LLMs to auto-complete these test cases, allowing developers to stay within their IDE. However, unlike NL-LLMs, whose annotation abilities across various NLP tasks are well studied (Zhang et al., 2023b; Tan et al., 2024; Uzair-UI-Haq et al., 2025), the effectiveness of Code-LLMs for this role remains underexplored, with previous work focusing on specific tasks (Li et al., 2023) and ad-hoc models (Mohajeri et al., 2024; Zhang et al., 2023a).

To address this gap, we propose CODE-PROMPT, a first comprehensive framework for evaluating Code-LLMs across various NLP tasks. When prompting Code-LLMs, we reformulate natural language prompts into code-aligned formats (see Figure 1) inspired by modern (Python) development practices. Each prompt consists of three components: (1) a function definition serving as the task description; (2) unit tests as demonstrations; and (3) a final, incomplete unit test for inference.

We evaluate CODE-PROMPT in zero-shot and

¹<https://github.com/features/copilot>

²<https://cursor.com/>

³<https://leetcode.com>

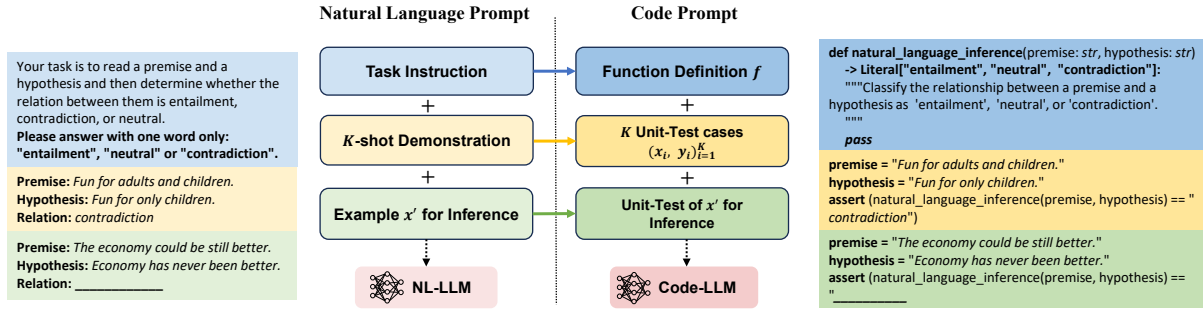


Figure 1: Overview of CODE-PROMPT, showcasing a natural language inference (NLI) task. On the left is the natural-language prompt sent to an NL-LLM, which is asked to generate a prediction after “Relation.” On the right is the corresponding code-style prompt for Code-LLMs, formulated as a Python unit-test suite over a defined NLI function. The Code-LLM completes the final assert statement as prediction.

few-shot settings and compare their performance against parallel NL-LLMs – models of the same size, version, and provider (see Section 3.2 and Table 3). We find that Code-LLMs using code prompts achieve comparable performance to NL-LLMs across multiple NLP tasks, while producing substantially fewer redundant tokens beyond the target label. Our cross-lingual experiments demonstrate strong knowledge transfer across different programming languages. However, a substantial performance gap remains across human languages, consistent with similar observations from NL-LLMs (Xuan et al., 2025). Given the rapid progress and increasing attention surrounding Code-LLMs (Yang et al., 2025a; Ge et al., 2025), we regard our work as an essential step in exploring their abilities beyond code.

2 Method

We now present CODE-PROMPT, showcased in Figure 1. We first introduce the prompt template and then describe the inference process. As our evaluation primarily uses prompts written in Python, we explain the framework using Python examples.

2.1 Prompt Construction

To leverage Code-LLMs for NLP tasks, we reformulate text generation under a natural language prompt (*NL prompt*) as a **code-completion** problem under a prompt written in code-style (*code prompt*). A code prompt typically contains three components: a function definition, unit-tests from labeled examples and an incomplete unit-test for the example to be predicted.

Function definition. We begin by defining a task-specific function f . The function name f_{name} reflects the underlying NLP task (Li et al., 2023) such

as `natural_language_inference`. The function arguments f_{args} correspond to the task input(s), such as `premise:str` and `hypothesis:str` for natural language inference (NLI).

Second, we convert the output specification from natural language instructions into an explicit Python type annotation f_{typing} . Take NLI for example, where the label set consists of *entailment*, *neutral*, and *contradiction*, we express f_{typing} as a `Literal` type enumerating these three labels, as illustrated in Figure 1. Then we put task description into a concise docstring f_{docstr} that clearly states the function’s purpose (Chen et al., 2021) (e.g., *classify the relationship between premise and hypothesis*) and reiterates the expected output (e.g., *as one of the target labels: ...*). Finally, the function body f_{body} is left empty (via `pass` in Python), since the function serves only as a semantic scaffold rather than an executable implementation. The final function definition is given by:

$$f := f_{\text{name}} \oplus f_{\text{args}} \oplus f_{\text{typing}} \oplus f_{\text{docstr}} \oplus f_{\text{body}}, \quad (1)$$

where \oplus denotes concatenation of code strings.

Unit-tests for demonstration. A unit test is expressed as an `assert` statement that applies the defined function f to an input x (mapped to variables from the function arguments f_{args}) and checks whether the output matches the available label y :

$$\text{assert } f(x) == y. \quad (2)$$

Inspired by prior work on few-shot demonstrations in natural-language prompting (Gao et al., 2021) and code-style prompting (Li et al., 2023), we convert each few-shot example into a corresponding `assert` statement and concatenate these statements to form the code-aligned demonstrations in the code prompt.

Category	Dataset	# Class	# Train	# Val	# Test	Classes
TC	SST-2 (Socher et al., 2013)	2	67,349	872	1821 [†]	positive, negative
	AG News (Zhang et al., 2015)	4	120,000	-	76,000	World, Sport, Business, Sci/Tech
	CoLA (Warstadt et al., 2019)	2	8551	1043	1063 [†]	acceptable, unacceptable
NLI	MRPC (Dolan and Brockett, 2005)	2	3668	408	1725	not_equivalent, equivalent
	XNLI (Conneau et al., 2018)	3	392,702	2490	5010	entailment, neutral, contradiction
	MSciNLI (Sadat and Caragea, 2024)	4	127,320	1000	4000	contrasting, reasoning, ...
RC	SemEval (Hendrickx et al., 2010)	10	8000	-	2717	Member-Collection, Cause-Effect, ...
	SciERC (Zhang et al., 2015)	7	3219	455	974	COMPARE, CONJUNCTION, ...
NR	Iris (Fisher, 1936)	3	150	-	-	setosa, versicolor, virginica
	HCC-Staging (Ours)	2	116	28	31	True (suitable for transplantation), False (unsuitable)

Table 1: Overview of the 10 benchmarked datasets. TC: Text classification. NLI: Natural language inference. RC: Relation classification. NR: Numeric reasoning. For XNLI, we report statistics for the English split, which are identical across the other 14 languages. †: The test split is not labelled.

Unit-test for inference. At inference, we construct a final `assert` statement using the test input x' and leave its expected output for the model to complete:

$$\text{assert } f(x') == \text{_____} \quad (3)$$

The Code-LLM is prompted to fill in the missing value on the right-hand side of the equality, which is interpreted as the model’s prediction.

Hence, the final code prompt is formed by concatenating (1) the function definition f , (2) the K unit-tests derived from the few-shot examples $\{(x_i, y_i)\}_{i=1}$, and (3) the incomplete unit test for the inference example x' . Appendix E presents the resulting code prompts for each individual task.

2.2 Inference

Code-LLMs first attempt to complete the final unit test for x' according to their underlying generation objective — normally autoregressively (Radford et al., 2019; Raffel et al., 2020), i.e. by iteratively selecting the most probable next token conditioned on all existing tokens (those provided in the prompt as well as those generated so far).

For classification tasks, the final prediction is determined by scanning the model’s output for the first occurrence of any valid target label. This post-processing is also applied to NL-LLM generations.

3 Experiments

The evaluated datasets (Section 3.1) and models (Section 3.2) are as follows, together with our evaluation setup (Section 3.3) including data scenarios and evaluation metrics. Further implementation details are listed in Appendix A.

3.1 Datasets

We collect 10 datasets from four categories of NLP tasks: text classification (TC), natural language inference (NLI), relation classification (RC) and numeric reasoning (NR). Table 1 provides an overview of all the benchmarked datasets.

3.1.1 Text classification

We evaluate CODE-PROMPT on sentiment classification (SST-2, Socher et al. (2013)), topic classification (AG News, Zhang et al. (2015)) and grammatical acceptability classification (CoLA, Warstadt et al. (2019)).

3.1.2 Natural language inference

NLI requires modeling a pair of texts to predict semantic relation between them (such as entailment). In this paper we benchmark three NLI datasets: XNLI (Conneau et al., 2018) to recognize entailment, contradict or neither (*neutral*) between premise and hypothesis sentences; MRPC (Dolan and Brockett, 2005) to determine if two texts construct paraphrases; and MSciNLI (Sadat and Caragea, 2024) for a new and challenging task of detecting more complicated relation from publications in Computer Science.

3.1.3 Relation classification

RC aims to recognize the relation of two extracted entities (Alt et al., 2020) in a text after a preliminary named entity recognition (Tjong Kim Sang and De Meulder, 2003) step. In this study, we consider SemEval (Hendrickx et al., 2010) in general domains and SciERC (Luan et al., 2018) in the scientific domain.

3.1.4 Logic-based numeric reasoning

Some NLP tasks require models to understand text but also apply explicit logics (Kazemi et al., 2023;

Input	Label	Explanation for labelling
35mm x 20mm x 15mm	<i>True</i>	Single tumor with diameter 35mm \leq 5 cm
2cm, 1cm, 15mm, 12mm	<i>False</i>	4 tumors in total, which exceeds the count limit 3
5x5mm, 17x28mm	<i>True</i>	2 tumors with maximal diameter 28mm \leq 3 cm

Table 2: Examples of annotation in HCC-Staging. The labeling follows Milan criteria introduced in Section 3.1.

Servantez et al., 2024); besides, the ability of understanding in text is also crucial (Yang et al., 2025b), especially in scientific context (Akhtar et al., 2023). Since Code-LLMs are pre-trained on code with programming logic and numbers, they might be well-suited for numeric and logical tasks. To evaluate the ability of logic-based numeric reasoning, we consider two datasets: **reformulated Iris** (Fisher, 1936) and a new clinical dataset **HCC-Staging** of patients with hepatocellular carcinoma (HCC).

The original Iris dataset classifies flowers into three species given their numeric features (i.e. sepal and petal dimensions). To reformulate this task into a benchmark for logic-based numeric reasoning, we relabel each example using the following logic:

If petal width < 0.8 cm, then label setosa;
If petal width \geq 0.8 cm and < 1.8 cm, then versicolor; *Otherwise, label virginica.*

The derived labels agree with the original with 0.96 (144/150) accuracy and serve as the reference labels in this paper. In reformulated Iris, each input is presented in a structured textual format (e.g., “Petal length: x cm, Petal width: y cm, ...”), requiring models to correctly identify the relevant measurement (petal width) and apply the logic-based numerical comparison to determine the label.

In the context of NLP-assisted clinical care, numeric reasoning is essential for interpreting quantitative information in medical reports to support treatment decisions (Borchert et al., 2022). In this study, we construct dataset **HCC-Staging** by collecting 175 HCC cases provided by an university hospital and extracting tumor sizes from them (see Appendix D for details). Each case is labeled according to *Milan criteria* (Mazzaferro et al., 1996), which determine a patient’s suitability for liver transplantation:

If the patient has one single tumor with diameter⁴ \leq 5 cm, then label as True (suitable for transplantation); Alternatively, if there are up

⁴Size of a tumor can be given in 1-D (e.g. 5 mm), 2-D (e.g. 2cm x 3cm) or 3-D (e.g. 3x4x5 mm). In multi-dimensional cases, diameter is defined as the largest dimension.

Provider	Model	# Params (B)
OpenAI	GPT-3.5 Turbo	175
DeepSeek	DeepSeek-V3.1	685
Qwen	Qwen3-Coder-30B	30.5
	Qwen3-30B-Instruct	30.5
	Qwen2.5-Coder-32B	32.8
	Qwen2.5-32B-Instruct	32.8
	Qwen2.5-Coder-14B	14.8
	Qwen2.5-14B-Instruct	14.8
	Qwen2.5-Coder-7B	7.6
	Qwen2.5-7B-Instruct	7.6
	Qwen2.5-Coder-3B	3.1
	Qwen2.5-3B-Instruct	3.1
Llama	CodeLlama-13b	13
	CodeLlama-13b-Python	13
	Llama-2-13b	13
	CodeLlama-7b	6.7
Gemma	CodeLlama-7b-Python	6.7
	Llama-2-7b	7
	codegemma-7b	8.5
	gemma-7b	8.5
	codegemma-2b	2.5
	gemma-2b	2.5

Table 3: Overview of the evaluated models.

to three tumors, and each has diameter \leq 3 cm, then also True; Otherwise, label False.

During inference, LLMs are given Milan criteria (via task instruction in NL prompt or docstring in code prompt) and required to apply them to classify each patient. Compared to Iris, this task requires abilities of not only numerical comparison but also counting and unit conversion (e.g., *mm* to *cm* if necessary). Table 2 illustrates example annotations.

3.2 Models

We aim to evaluate **paired** Code-LLMs and NL-LLMs — i.e., models that come from the same provider with comparable release dates and parameter scales — so that differences in performance can be attributed to pre-training and prompting style rather than unrelated architectural disparities. Hence, our model selection follows two criteria: (1) we identify publicly available Code-LLMs (as of 15.09.2025), either open-weight or accessible via a public code completion API⁵, from five major LLM providers: OpenAI (Brown et al., 2020; OpenAI et al., 2024), DeepSeek (Guo et al., 2024), Qwen (Hui et al., 2024), Llama (Rozière et al., 2023) and Gemma (Zhao et al., 2024). Notably, **GPT-3.5 Turbo** is the last OpenAI model that provides a public complete API capable of code completion; newer GPT models are restricted to chat-

⁵For example, <https://api.openai.com/v1/completions> from OpenAI and <https://api.deepseek.com/beta> from DeepSeek.

based APIs designed primarily for NL interaction. (2) For each identified Code-LLM, we include it only if a parallel NL-LLM counterpart exists, defined as a model from the same provider with comparable version (or release date) and parameter size. Pairs for which one side (Code or NL) is unavailable are excluded to ensure controlled comparison.

After the two steps, we include a total of 22 large language models (grouped in Table 3). Notably, GPT-3.5 Turbo and DeepSeek-V3.1 can be regarded as both NL- and Code-LLMs since they provide a unified complete API capable of both natural language and code generation. For these two models, we use identical backbones but different prompt formats, referring to them as X (Code) and X (NL), respectively. In this paper, to ensure fair and reproducible comparisons, we evaluate each model in full precision (no quantization) and generate with a fixed temperature of 0.

3.3 Evaluation Setup

Data scenarios. We evaluate CODE-PROMPT in zero-shot and few-shot scenarios.

Zero-shot requires the model to predict based on solely prompt instructions without any demonstrated example. For most datasets, we evaluate on the full test split. SST-2 and CoLA are exceptions because their test splits are unlabelled; for these two tasks, we instead evaluate on the validation split. For Iris, which does not provide any data division, we evaluate on the entire dataset.

Few-shot setting is often cast as an N -way K -shot problem (Vinyals et al., 2016), where K examples from each of the N classes are available as demonstration (Li et al., 2023). For datasets with predefined splits, we randomly sample $N \times K$ examples from the training split and evaluate on the same test (or validation) data used in the zero-shot setting. The only exception is Iris, which has no predefined splits; in this case, we sample $N \times K$ examples from the full dataset for prompting and evaluate on all remaining (unsampled) instances.

The selection of K is constrained by each model’s context window, the number of classes, and the length of the demonstrations. For most datasets we set K to be 4, which fits within the context limits of all models. Relation classification datasets contain substantially more classes, making larger K infeasible: for SCIERC we set K as 2, and for SemEval even a 1-shot prompt exceeds the context window of several models. Consequently, we exclude SemEval from the few-shot experiments.

To ensure fair comparison and reproducible results, we repeat the few-shot evaluation three times using fixed random seeds 0, 1 and 2. Therefore, within each run the same set of sampled shots is used across all test instances and all models.

Evaluation metrics. For each dataset, we report the performance scores commonly applied in prior research. Following Wang et al. (2018), we report accuracy for SST-2 and MRPC, and Matthews correlation for CoLA; following Zhang et al. (2015), we report accuracy for AG News; following Conneau et al. (2018), we report accuracy for XNLI; following Sadat and Caragea (2024), we report Macro-F1 for MSciNLI; we report Micro-F1 for SemEval and SCIERC according to Harbecke et al. (2022) and Luan et al. (2018), respectively; we report accuracy for Iris following Nugroho et al. (2020) and apply the same metric for HCC-Staging.

Beyond correctness, we introduce redundancy, measuring how well models adhere to the “answer only with the label” instruction in both NL and code prompts (see Appendix E). Redundancy R is defined as the complement of the ratio between the label length and the total generation:

$$R := 1 - \frac{\text{length}(\text{label})}{\text{length}(\text{generation})} \quad (4)$$

where length is measured in characters rather than tokens to avoid tokenizer-specific variation across models. Special tokens such as `<eos>` are excluded in computation. A lower redundancy indicates that the model produces fewer extratokens and adheres more strictly to the prompt. In this study we report both performance (implemented by sklearn (Pedregosa et al., 2011)) and redundancy.

4 Results and Analysis

We first present results in zero-shot and few-shot scenarios in general; then we study the effect of multilinguality, including different programming languages as well as human languages.

4.1 Zero-shot Results

Table 4 presents the zero-shot performance for different models and tasks. We see that Code-LLMs deliver comparable performance with NL-LLMs on all NLP tasks: (1) Code-LLMs outperform NL-LLMs in 7 of the 11 pair-wise comparisons, measured by average accuracy \bar{X} over all ten datasets. In the remaining cases, their largest deficit to NL model is 6.5% (from Qwen2.5-14B). (2) The

Model	TC			NLI			RC		NR		\bar{X}
	SST-2	AG News	CoLA	MRPC	XNLI	MSciNLI	SemEval	SciERC	Iris	HCC-Staging	
	Acc	Acc	MCC	Acc	Acc	Macro-F1	Micro-F1	Micro-F1	Acc	Acc	
GPT-3.5 Turbo (Coder)	94.6	61.3	57.6	73.7	68.0	37.5	45.2	44.6	92.7	58.1	63.3
GPT-3.5 Turbo (NL)	93.1	62.5[†]	64.1	76.7	68.4	34.2	48.9	50.4	96.0	64.5	65.9
DeepSeek-V3.1 (Code)	95.1	60.0	66.6	76.6	75.6	60.4	62.0[†]	59.2	96.7[†]	83.9[†]	73.6[†]
DeepSeek-V3.1 (NL)	94.6	56.7	67.9[†]	78.5	83.3	63.1[†]	29.0	34.0	96.0	80.6	68.4
Qwen3-Coder-30B	94.0	60.8	60.9	79.2[†]	84.2	49.8	46.8	45.9	96.0	67.7	68.5
Qwen3-30B-Instruct	95.3[†]	61.7	67.2	74.7	73.8	41.2	49.1	53.6	30.0	74.2	62.1
Qwen2.5-Coder-32B	94.8	50.4	49.7	77.7	77.8	44.0	42.6	59.8[†]	94.5	80.6	67.2
Qwen2.5-32B-Instruct	93.5	61.4	67.1	78.5	89.1[†]	58.2	57.3	49.0	96.0	64.5	71.5
Qwen2.5-Coder-14B	94.5	51.6	46.1	73.4	82.2	46.4	33.1	56.8	68.7	71.0	62.4
Qwen2.5-14B-Instruct	89.8	59.2	65.1	77.0	84.6	52.6	51.2	57.8	90.0	61.3	68.9
Qwen2.5-Coder-7B	93.8	49.1	49.6	69.4	77.4	44.4	31.3	18.4	86.0	48.4	56.8
Qwen2.5-7B-Instruct	89.6	57.8	58.5	50.6	80.3	22.8	42.7	46.0	78.7	67.7	59.5
Qwen2.5-Coder-3B	90.1	32.5	40.0	69.7	65.8	42.1	31.8	49.8	72.7	51.6	54.6
Qwen2.5-3B-Instruct	88.9	52.9	52.8	49.9	75.9	22.0	36.5	47.8	33.3	64.5	52.5
CodeLlama-13b	89.0	37.9	30.3	68.0	57.6	22.9	36.3	14.9	54.0	51.6	46.3
CodeLlama-13b-Python	88.1	48.7	22.0	74.0	56.1	35.2	29.6	14.9	34.0	51.6	45.4
Llama-2-13b	39.0	14.7	3.7	27.8	44.3	19.5	20.4	15.2	0.0	48.4	23.3
CodeLlama-7b	89.3	53.3	14.7	66.7	38.9	15.6	23.9	14.0	37.3	48.4	40.2
CodeLlama-7b-Python	89.3	54.9	26.4	69.1	53.4	20.2	25.3	17.1	33.3	48.4	43.7
Llama-2-7b	42.1	14.5	-6.0	22.6	20.0	11.9	20.5	2.4	0.0	35.5	16.4
codegemma-7b	91.9	51.0	17.6	73.6	57.2	28.0	35.3	30.9	66.7	67.7	52.0
gemma-7b	41.2	18.9	5.4	63.7	43.8	24.7	21.6	7.7	66.7	48.4	34.2
codegemma-2b	79.4	31.4	0.0	66.5	34.4	10.4	25.6	12.9	33.3	48.4	34.2
gemma-2b	51.6	24.0	2.9	31.6	19.5	11.6	20.3	21.1	33.3	45.2	26.1
Avg. Code-LLM	91.1	49.5	37.0	72.1	63.7	35.1	36.1	33.8	66.6	59.8	54.5
Avg. NL-LLM	74.4	44.0	40.8	57.4	62.1	32.9	36.1	35.0	56.4	59.5	49.9
Avg. NL-LLM (w/o Llama, gemma)	92.1	58.9	63.2	69.4	79.3	42.0	45.0	48.4	74.3	68.2	64.1

Table 4: Zero-shot performance (%) of NL-LLMs using natural language prompts and Code-LLMs using code-based prompts. The higher score within each comparison is highlighted in bold. The highest score of each task is denoted with [†]. *Acc*: Accuracy. *MCC*: Matthews correlation coefficient. \bar{X} denotes the arithmetic mean across 10 tasks.

Model	TC			NLI			RC		NR		\bar{R}
	SST-2	AG News	CoLA	MRPC	XNLI	MSciNLI	SemEval	SciERC	Iris	HCC-Staging	
GPT-3.5 Turbo (Code)	0.0	26.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2.6
GPT-3.5 Turbo (NL)	0.0	25.2	0.0	0.0	0.0	0.0	7.6	10.7	0.0	0.0	4.4
DeepSeek-V3.1 (Code)	0.0	24.4	0.0	0.0	0.2	0.0	0.0	0.0	0.0	0.0	2.5
DeepSeek-V3.1 (NL)	0.0	25.3	0.0	0.0	0.0	0.0	92.5	77.3	0.0	0.0	19.5
Qwen3-Coder-30B	20.0	36.7	16.1	13.7	16.3	19.0	6.9	10.3	19.9	30.9	19.0
Qwen3-30B-Instruct	88.5	93.0	85.2	84.1	83.8	89.6	79.0	87.4	97.4	90.8	87.9
Qwen2.5-Coder-32B	0.0	29.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2.9
Qwen2.5-32B-Instruct	71.6	83.5	81.4	86.8	88.3	88.5	78.8	87.3	51.0	88.6	80.6
Qwen2.5-Coder-14B	0.0	26.1	0.0	0.0	0.0	0.0	0.0	0.1	0.0	0.0	2.6
Qwen2.5-14B-Instruct	80.7	92.1	83.5	86.9	88.6	88.8	78.5	87.8	77.3	92.2	85.6
Qwen2.5-Coder-7B	0.0	21.4	0.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2.2
Qwen2.5-7B-Instruct	89.2	92.7	86.8	85.7	83.6	90.0	80.2	88.4	89.3	92.0	87.8
Qwen2.5-Coder-3B	0.0	5.9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.6
Qwen2.5-3B-Instruct	89.2	93.1	85.5	68.1	86.5	91.1	80.1	88.6	94.9	93.2	87.0
CodeLlama-13b	84.6	89.9	80.4	80.0	78.3	79.3	71.4	81.5	83.4	91.8	82.1
CodeLlama-13b-Python	84.7	89.2	77.1	78.0	78.3	79.7	77.5	74.5	87.1	91.3	81.7
Llama-2-13b	79.6	94.0	75.9	99.3	97.3	95.5	52.5	13.0	100.0	90.8	79.8
CodeLlama-7b	84.1	85.8	79.2	79.2	75.9	78.3	59.4	77.9	81.5	89.5	79.1
CodeLlama-7b-Python	11.1	82.4	20.7	74.6	31.4	78.1	62.5	70.3	85.9	93.6	61.1
Llama-2-7b	80.4	93.6	96.6	98.8	99.8	99.6	46.4	89.0	100.0	94.3	89.9
codegemma-7b	0.0	24.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2.4
gemma-7b	50.5	92.1	87.1	83.6	77.4	77.6	79.8	89.5	94.4	92.4	82.4
codegemma-2b	0.0	0.9	67.9	84.9	4.6	0.5	0.1	0.7	0.0	0.0	16.0
gemma-2b	81.8	89.4	90.5	91.2	91.1	70.2	81.0	75.3	0.0	3.2	67.4
Avg. Code-LLM	21.9	41.7	26.3	31.6	21.9	25.8	21.4	24.3	27.5	30.5	27.3
Avg. NL-LLM	64.7	79.5	70.2	71.3	72.4	71.9	68.8	72.2	64.0	67.0	70.2

Table 5: Redundancy (%) (defined in Section 3.3; the lower, the better) of zero-shot experiments. The lower redundancy in each comparison is highlighted in bold. \bar{R} denotes mean redundancy across 10 tasks.

mean accuracy across all Code-LLMs is 54.5%, exceeding that of NL-LLMs (49.9%). Excluding four underperforming NL-LLMs (from Llama and Gemma) as outliers, the remaining NL-LLMs achieve a mean accuracy of 64.1%, outperforming Code-LLMs by 9.6%. (3) On a per dataset basis, NL-LLMs achieve the best results (marked with †) on all three text classification and two out of the three NLI datasets. In contrast, Code-LLMs achieve the top performance on both numerical reasoning (NR) tasks. We attribute this divergence to task characteristics, since TC and NLI rely pri-

marily on NL understanding and NR requires more understanding of numbers and logics.

Table 5 shows the **redundancy** for different models. We observe remarkably lower redundancy from Code-LLMs: (1) in 10 out of 11 pair-wise comparisons, code models achieve better average redundancy across tasks (indicated by \bar{R}) than NL models. (2) The mean redundancy across all Code-LLMs is 27.3%, markedly better than 70.2% from NL-LLMs. (3) Across all experiments in Table 5, Code-LLMs produce zero-redundancy outputs in 63 runs (48.5% of all their runs), whereas NL-

# Shots	SST-2 4	AG News 4	CoLA 4	MRPC 4	XNLI 4	MSciNLI 4	SciERC 2	Iris 4	HCC-Staging 4	\bar{X}
Performance	<i>Acc</i>	<i>Acc</i>	<i>MCC</i>	<i>Acc</i>	<i>Acc</i>	<i>Macro-F1</i>	<i>Micro-F1</i>	<i>Acc</i>	<i>Acc</i>	\bar{X}
Avg. Code-LLM	91.7	54.7	42.7	72.0	70.1	41.3	30.8	87.4	69.7	62.2
Avg. NL-LLM	59.9	39.4	39.0	58.3	65.2	36.8	42.7	62.5	51.3	50.6
Avg. NL-LLM (w/o LLaMA, Gemma)	93.1	61.1	60.9	63.7	84.2	47.3	59.2	93.5	71.0	70.4
Δ Performance compared to zero-shot										Δ
Avg. Code-LLM	+0.6	+5.3	+5.6	-0.2	+6.3	+6.1	-3.0	+20.8	+9.9	+5.7
Avg. NL-LLM	-14.5	-4.7	-1.7	+0.9	+3.1	+3.9	+7.7	+6.1	-8.2	-0.8
Avg. NL-LLM (w/o LLaMA, Gemma)	+1.0	+2.2	-2.4	-5.7	+4.8	+5.2	+10.9	+19.2	+2.8	+4.2
Redundancy (the lower, the better)										\bar{R}
Avg. Code-LLM	17.9	34.8	23.3	19.5	29.3	17.6	24.9	22.0	20.2	23.3
Avg. NL-LLM	64.6	82.6	73.2	64.7	71.2	75.7	69.7	55.1	63.4	68.9
Avg. NL-LLM (w/o LLaMA, Gemma)	58.9	73.0	60.3	54.7	62.0	63.9	59.3	60.3	65.4	62.0
Δ Redundancy compared to zero-shot (minus value means improvement)										Δ
Avg. Code-LLM	-3.9	-6.9	-3.0	-12.1	+7.4	-8.2	+0.7	-5.6	-10.4	-4.7
Avg. NL-LLM	-0.1	+3.2	+3.0	-6.7	-1.2	+3.8	-2.5	-9.0	-3.7	-1.5
Avg. NL-LLM (w/o LLaMA, Gemma)	-1.0	+0.9	0.0	-4.1	+0.5	-0.1	-16.0	+1.7	+0.1	-2.0

Table 6: Aggregated performance and redundancy (%) of few-shot experiments, with their change (Δ) from zero-shot highlighted in bold. Full few-shot results of all models are presented in Table 12 of performance and Table 13 of redundancy. SemEval results are missing due to a limited context window for this task (see Section 3.3).

Model	SST-2 (Acc)				MRPC (Acc)				SemEval (Micro-F1)				HCC-Staging (Acc)				\bar{X}
	Python	JS	C++	Ruby	Python	JS	C++	Ruby	Python	JS	C++	Ruby	Python	JS	C++	Ruby	
GPT-3.5 Turbo (Code)	94.6	94.4	94.5	95.1	73.7	72.1	74.6	71.6	45.2	46.0	44.6	46.2	58.1	58.1	64.5	58.1	62.4
DeepSeek-V3.1 (Code)	95.1	95.2	95.2	95.6	76.6	78.9	77.0	77.0	62.0	65.1	62.3	62.5	83.9	77.4	87.1	77.4	75.6
Qwen3-Coder-30B	94.0	95.0	94.5	94.2	79.2	79.1	79.4	77.3	46.8	41.8	45.7	43.6	67.7	74.2	64.5	64.5	62.1
CodeLlama-13b	89.0	90.6	90.5	90.4	68.0	68.2	69.7	69.3	36.3	39.9	36.5	36.0	51.6	51.6	51.6	51.6	55.0
codegemma-7b	91.9	92.2	92.5	92.2	66.7	68.1	67.7	67.8	35.3	41.8	41.8	37.7	67.7	71.0	67.7	67.7	63.4
Avg.	92.5	93.3	93.2	93.1	72.6	73.6	73.5	72.9	45.1	47.2	46.6	45.0	67.7	68.6	67.7	65.3	

Table 7: Zero-shot performance (%) of four different programming languages: Python, JavaScript (JS), C++ and Ruby. Code-prompts written in those programming languages are presented in Appendix E.

LLMs do so in only 15 (13.6%). We further notice even smaller, open-source coders such as Qwen-Coder and codegemma often output only the target label, adhering strictly to the format specification. By contrast, most evaluated open-source NL-LLMs struggle with following concise output constraints and frequently add unsolicited text (e.g. *Based on the information or I am not sure*). The findings indicate an inherent strength of Code-LLMs that they are naturally better at generating precise, format-faithful outputs.

4.2 Few-shot Results

We present few-shot results (performance in Table 12 and redundancy in Table 13) in Appendix B. We compare them with zero-shot results based on average scores across all NL-/code-models (see Table 6). We observe that Code-LLMs benefit more consistently from few-shot demonstrations: (1) Code-LLMs show performance gains on 7 of 9 tasks, whereas NL-LLMs improve on only 5. (2) Averaged across all tasks, Code-LLMs achieve a 5.7% improvement, while NL-LLMs exhibit a slight overall decline (-0.8%). Since we notice this decline is driven by four underperforming NL-LLMs (from Llama and gemma), we exclude them

as outliers; the remaining NL-LLMs show a more moderate improvement of 4.2%.

A similar trend appears in **redundancy**. Across all tasks, Code-LLMs reduce redundancy by 4.7 percentage points, compared to 1.5 for NL-LLMs (or 0.2 for the refined NL-LLM subset mentioned above). These results indicate that Code-LLMs not only learn more effectively from few-shot demonstrations but also produce more concise and format-compliant outputs under few-shot prompting.

4.3 Results of Multilinguality

To assess multilingual abilities of Code-LLMs, we first select the best performing (indicated by highest mean score across tasks \bar{X} in Table 4) code model from each of the 5 providers: GPT-3.5 Turbo, DeepSeek-V3.1, Qwen3-Coder-30B, CodeLlama-13b and codegemma-7b; then we experiment in different programming and human languages.

Programming languages. We translate code prompts (originally in Python) into three additional programming languages: two high-resource (JavaScript and C++) and one relatively low-resourced (Ruby). All four languages share a property: they support one-line, assert-style unit tests, enabling a consistent prompt structure across lan-

Model	AR	BG	DE	EL	EN	ES	FR	HI	RU	SW	TH	TR	UR	VI	ZH	Avg.
GPT-3.5 Turbo (Code)	57.8	60.4	63.5	59.0	68.0	63.3	63.2	52.2	59.2	56.3	52.6	58.0	49.2	58.7	59.9	58.8
DeepSeek-V3.1 (Code)	71.0	74.2	74.5	73.4	75.6	74.4	75.8	70.6	73.4	64.5	70.1	71.5	66.1	73.0	74.3	72.2
Qwen3-Coder-30B	72.5	75.3	75.1	73.2	84.2	79.1	77.7	68.1	74.0	62.4	69.4	71.6	61.2	74.5	76.0	73.0
CodeLlama-13b	45.7	51.5	52.5	43.4	57.6	54.1	54.5	40.1	53.4	33.7	43.2	47.8	42.5	48.3	49.2	47.8
codegemma-7b	49.0	51.9	53.8	49.5	57.2	55.8	53.9	45.0	52.5	46.2	49.0	51.7	42.2	49.5	50.4	50.5
Avg.	59.2	62.7	63.9	59.7	68.5	65.3	65.0	55.2	62.5	52.6	56.9	60.1	52.2	60.8	62.0	
Ranking by Avg.	#11	#5	#4	#10	#1	#2	#3	#13	#7	#14	#12	#9	#15	#8	#6	

Table 8: Zero-shot accuracy (%) of XNLI over 15 languages. The prompt is presented in Appendix E.5. Evaluated languages are AR (Arabic), BG (Bulgarian), DE (German), EL (Greek), EN (English), ES (Spanish), FR (French), HI (Hindi), RU (Russian), SW (Swahili), TH (Thai), TR (Turkish), UR (Urdu), VI (Vietnamese), and ZH (Chinese).

456 guages (see Appendix E for all translated prompts).

457 We conduct zero-shot experiments with trans- 497
458 lated code prompts on four datasets, one drawn 498
459 from each task category. As shown in Table 7, 499
460 performance remains highly consistent across pro- 500
461 gramming languages. For each task, we average 501
462 scores over five models per language and then com- 502
463 pare these language-wise averages. The resulting 503
464 divergences are minimal: 0.8% for SST, 1.0% for 504
465 MRPC, 2.1% for SemEval, and 3.3% for HCC- 505
466 Staging — indicating that the world knowledge 506
467 encoded in Code-LLMs to perform evaluated NLP 507
468 tasks is largely programming-language agnostic, 508
469 transferring robustly across these 4 languages. 509

470 **Human languages.** We further evaluate Code- 510
471 LLMs on non-English inputs using dataset XNLI. 511
472 Since XNLI provides sentence-aligned translations 512
473 across 15 languages, it allows direct, language- 513
474 wise comparison. We keep the prompt template in 514
475 English and add a short instruction indicating the 515
476 language such as “Both premise and hypothesis are 516
477 written in German”. Table 8 shows Code-LLMs 517
478 perform best on English (68.5%), followed by other 518
479 high-resource European languages such as Spanish, 519
480 French and German. A substantial drop appears for 520
481 lower-resource languages (e.g., -16.3% for Urdu, 521
482 -15.9% for Swahili). Table 9 in Appendix B spot- 522
483 lights these five Code-LLMs across English and 523
484 Chinese. The two China-based models (DeepSeek, 524
485 Qwen-Coder) exhibit the smallest relative accuracy 525
486 drop when switching from English to Chinese. 526

487 5 Related Work

488 **Evaluating Code-LLMs** has largely focused on 527
489 code generation, either from natural language in- 528
490 structions (e.g., MBPP (Austin et al., 2021), CODE- 529
491 JUDGE (Tong and Zhang, 2024), CodeArena (Yang 530
492 et al., 2025c)) or from partial code contexts (*auto-* 531
493 *completion*), such as HumanEval (Chen et al., 532
494 2021) and RepoBench (Liu et al., 2024). Some 533
495 recent benchmarks extend this scope to broader 534
535

496 software engineering tasks, including unit test gen- 496
497 eration (Yang et al., 2024) and Github issue-driven 497
498 code modification (Jimenez et al., 2024). To our 498
499 knowledge, the general NLP capabilities of Code- 499
500 LLMs remain largely unexplored. 500

501 However, **applying Code-LLMs to NLP** has 501
502 shown promise in specific settings. Li et al. (2023) 502
503 show that code-style prompting with Code-LLMs 503
504 can outperform NL-LLMs on information extrac- 504
505 tion tasks in few-shot scenarios and introduce *for-* 505
506 *mat fidelity* to assess whether generated code con- 506
507 forms to a parsable Python dictionary structure. 507
508 Mohajeri et al. (2024) evaluate CodeLlama on 508
509 three TC and one NLI task, also under few-shot 509
510 settings. Zhang et al. (2023a) explore code-style 510
511 prompts with GPT models for generation tasks such 511
512 as QA and summarization, again relying on few- 512
513 shot demonstrations. Our work extends this line of 513
514 research through a systematic evaluation of Code- 514
515 LLMs over a broad spectrum of NLP tasks and 515
516 model families, across both zero-shot and few-shot 516
517 regimes and multilingual settings. 517

518 6 Conclusion

519 In this paper, we present CODE-PROMPT, a uni- 519
520 fied evaluation framework for benchmarking Code- 520
521 LLMs on a diverse set of NLP tasks using code- 521
522 style prompts. The framework is readily extensible 522
523 beyond the tasks considered in this study. We evalu- 523
524 ate 22 LLMs (including 13 Code-LLMs) across 524
525 10 datasets under zero-shot and few-shot settings, 525
526 and further investigate the impact of different pro- 526
527 gramming and human languages. Our results show 527
528 that Code-LLMs achieve performance comparable 528
529 to NL-LLMs on a wide range of NLP tasks, while 529
530 exhibiting distinctive traits such as reduced redun- 530
531 dant generation and strong transferability across 531
532 programming languages. We hope this work sheds 532
533 light on the natural language capabilities of Code- 533
534 LLMs and motivates further research at the inter- 534
535 section of code intelligence and NLP. 535

536 Limitations

537 **Potential data leakage in LLMs.** A major
538 source of bias in benchmarking large language
539 models is data leakage (Xu et al., 2024; Zhou
540 et al., 2025), where a benchmarked dataset may
541 have been included in a model’s pre-training corpus.
542 Such leakage can inflate reported performance
543 and undermine fair comparison. While many of
544 the datasets evaluated in this work were released
545 prior to the deployment of modern LLMs, Table 1
546 shows that MSciNLI (released in 2024) and our
547 proposed HCC-Staging dataset are relatively recent,
548 and therefore less likely to be present in
549 model pre-training data. We find that our main
550 conclusions remain consistent when focusing on
551 these two datasets: Code-LLMs achieve performance
552 comparable to NL-LLMs, exhibit substantially
553 better redundancy, and (on HCC-Staging)
554 demonstrate strong transferability across programming
555 languages. Nevertheless, we cannot fully rule
556 out indirect exposure through continuously updated
557 training pipelines, and the results should therefore
558 be interpreted with this limitation in mind.

559 **Limited access to the latest models.** As discussed
560 in Section 3.2, our evaluation includes the most recent
561 Code-LLMs that are publicly accessible at the time of
562 study, either as open-weight checkpoints or via a public
563 code-completion API. However, evaluating the very latest
564 coding models is constrained by limited openness in the
565 ecosystem: (1) several proprietary model families
566 that power state-of-the-art coding assistants (e.g.,
567 Gemini and Claude) are closed-source and do not
568 provide a public code-completion API. (2) Some
569 open-weight model lines, such as CodeLlama and
570 codegemma, have not released updated code-specialized
571 variants beyond the versions evaluated in this work.

572 We view this limitation as highlighting the importance
573 of **continued openness** in (Code-)LLM development
574 for future research, and we acknowledge providers
575 such as DeepSeek and Qwen for consistently releasing
576 up-to-date coding models to the research community.
577
578
579

580 Ethics Statement

581 The HCC-Staging dataset is derived from clinical
582 reports and therefore involves sensitive patient
583 information. To ensure ethical compliance and protect
584 patient privacy, we take the following measures: (1)
585 The source hospital anonymizes the original

586 reports by removing direct identifiers such as
587 patient names and dates of birth (replaced by age).
588 (2) We will publish the ethic code from our partner
589 hospital upon acceptance. (3) To further reduce
590 re-identification risk, we slightly perturbed all tumor
591 dimensions (e.g., 17 mm adjusted to 15 mm, 9 x 9 x
592 8 mm adjusted to 8 x 8 x 5 mm) while preserving
593 the original staging outcome under the Milan criteria.
594 The final post-processed dataset was reviewed
595 and approved by the annotator. Here we explicitly
596 specify non-commercial research-only use for this
597 HCC-Staging dataset.

References 598

- 599 Mubashara Akhtar, Abhilash Shankarampeta, Vivek
600 Gupta, Arpit Patil, Oana Cocarascu, and Elena Simperl.
601 2023. [Exploring the numerical reasoning capabilities of language models: A comprehensive analysis on tabular data](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages
602 15391–15405, Singapore. Association for Computational
603 Linguistics. 604
605
606
607 Christoph Alt, Aleksandra Gabryszak, and Leonhard
608 Hennig. 2020. [Probing linguistic features of sentence-level representations in neural relation extraction](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*,
609 pages 1534–1545, Online. Association for Computational
610 Linguistics. 611
612
613
614 Jacob Austin, Augustus Odena, Maxwell I. Nye,
615 Maarten Bosma, Henryk Michalewski, David Dohan,
616 Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le,
617 and Charles Sutton. 2021. [Program synthesis with large language models](#). *CoRR*, abs/2108.07732. 618
619
620 Florian Borchert, Christina Lohr, Luise Modersohn,
621 Jonas Witt, Thomas Langer, Markus Follmann,
622 Matthias Gietzelt, Bert Arnrich, Udo Hahn, and
623 Matthieu-P. Schapranow. 2022. [GGPONC 2.0 - the German clinical guideline corpus for oncology: Curation workflow, annotation policy, baseline NER taggers](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3650–
624 3660, Marseille, France. European Language Resources
625 Association. 626
627
628
629 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
630 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
631 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
632 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
633 Gretchen Krueger, Tom Henighan, Rewon Child,
634 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
635 Winter, and 12 others. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901.
636 Curran Associates, Inc. 637
638
639 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
640 Henrique Pondé de Oliveira Pinto, Jared Kaplan,

641	Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. Evaluating large language models trained on code . <i>CoRR</i> , abs/2107.03374.	697
642		698
643		699
644		700
645		701
646		702
647	Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: Evaluating cross-lingual sentence representations . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2475–2485, Brussels, Belgium. Association for Computational Linguistics.	703
648		704
649		705
650		706
651		707
652		708
653		709
654		710
655	William B. Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases . In <i>Proceedings of the Third International Workshop on Paraphrasing (IWP2005)</i> .	711
656		712
657		713
658		714
659	Kaiyue Feng, Yilun Zhao, Yixin Liu, Tianyu Yang, Chen Zhao, John Sous, and Arman Cohan. 2025. Physics: Benchmarking foundation models on university-level physics problem solving . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 11717–11743, Vienna, Austria. Association for Computational Linguistics.	715
660		716
661		717
662		718
663		719
664		720
665		721
666		722
667	Rory A. Fisher. 1936. The use of multiple measurements in taxonomic problems . <i>Annals of Human Genetics</i> , 7:179–188.	723
668		724
669	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3816–3830, Online. Association for Computational Linguistics.	725
670		726
671		727
672		728
673		729
674		730
675		731
676		732
677	Yuyao Ge, Lingrui Mei, Zenghao Duan, Tianhao Li, Yujia Zheng, Yiwei Wang, Lexin Wang, Jiayu Yao, Tianyu Liu, Yujun Cai, Baolong Bi, Fangda Guo, Jiafeng Guo, Shenghua Liu, and Xueqi Cheng. 2025. A survey of vibe coding with large language models . <i>CoRR</i> , abs/2510.12399.	733
678		734
679		735
680		736
681		737
682		738
683	Mingyang Geng, Shangwen Wang, Dezun Dong, Hao-tian Wang, Ge Li, Zhi Jin, Xiaoguang Mao, and Xiangke Liao. 2024. Large language models are few-shot summarizers: Multi-intent comment generation via in-context learning . In <i>Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, ICSE 2024, Lisbon, Portugal, April 14-20, 2024</i> , pages 39:1–39:13. ACM.	739
684		740
685		741
686		742
687		743
688		744
689		745
690		746
691	Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. Deepseek-coder: When the large language model meets programming - the rise of code intelligence . <i>CoRR</i> , abs/2401.14196.	747
692		748
693		749
694		750
695		751
696		752
	David Harbecke, Yuxuan Chen, Leonhard Hennig, and Christoph Alt. 2022. Why only micro-f1? class weighting of measures for relation classification . In <i>Proceedings of NLP Power! The First Workshop on Efficient Benchmarking in NLP</i> , pages 32–41, Dublin, Ireland. Association for Computational Linguistics.	753
		754
	Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2010. SemEval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals . In <i>Proceedings of the 5th International Workshop on Semantic Evaluation</i> , pages 33–38, Uppsala, Sweden. Association for Computational Linguistics.	755
		756
	Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, An Yang, Rui Men, Fei Huang, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. 2024. Qwen2.5-coder technical report . <i>CoRR</i> , abs/2409.12186.	757
		758
	Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: Can language models resolve real-world github issues? In <i>The Twelfth International Conference on Learning Representations</i> .	759
		760
	Mehran Kazemi, Quan Yuan, Deepti Bhatia, Najoung Kim, Xin Xu, Vaiva Imbrasaite, and Deepak Ramachandran. 2023. Boardgameqa: a dataset for natural language reasoning with contradictory information . In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23</i> , Red Hook, NY, USA. Curran Associates Inc.	761
		762
	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention . In <i>Proceedings of the 29th Symposium on Operating Systems Principles, SOSP '23</i> , page 611–626, New York, NY, USA. Association for Computing Machinery.	763
		764
	Peng Li, Tianxiang Sun, Qiong Tang, Hang Yan, Yuanbin Wu, Xuanjing Huang, and Xipeng Qiu. 2023. CodeIE: Large code generation models are better few-shot information extractors . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15339–15353, Toronto, Canada. Association for Computational Linguistics.	765
		766
	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing . <i>ACM Comput. Surv.</i> , 55(9).	767
		768
	Tianyang Liu, Canwen Xu, and Julian J. McAuley. 2024. Repubench: Benchmarking repository-level	769
		770

755	code auto-completion systems . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
756		
757		
758		
759	Do Xuan Long, Duy Dinh, Ngoc-Hai Nguyen, Kenji Kawaguchi, Nancy F. Chen, Shafiq Joty, and Min-Yen Kan. 2025. What makes a good natural language prompt? In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5835–5873, Vienna, Austria. Association for Computational Linguistics.	
760		
761		
762		
763		
764		
765		
766	Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3219–3232, Brussels, Belgium. Association for Computational Linguistics.	
767		
768		
769		
770		
771		
772		
773	Qinyu Luo, Yining Ye, Shihao Liang, Zhong Zhang, Yujia Qin, Yaxi Lu, Yesai Wu, Xin Cong, Yankai Lin, Yingli Zhang, Xiaoyin Che, Zhiyuan Liu, and Maosong Sun. 2024. RepoAgent: An LLM-powered open-source framework for repository-level code documentation generation . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 436–464, Miami, Florida, USA. Association for Computational Linguistics.	
774		
775		
776		
777		
778		
779		
780		
781		
782		
783	V Mazzaferro, E Regalia, R Doci, S Andreola, A Pulvirenti, F Bozzetti, F Montalto, M Ammatuna, A Morabito, and L Gennari. 1996. Liver transplantation for the treatment of small hepatocellular carcinomas in patients with cirrhosis. <i>N. Engl. J. Med.</i> , 334(11):693–699.	
784		
785		
786		
787		
788		
789	Mohammad Mahdi Mohajeri, Mohammad Javad Dousti, and Majid Nili Ahmadabadi. 2024. Cocop: Enhancing text classification with LLM through code completion prompt . <i>CoRR</i> , abs/2411.08979.	
790		
791		
792		
793	Heru Nugroho, Nugraha Priya Utama, and Kridanto Surendro. 2020. Performance evaluation for class center-based missing data imputation algorithm . In <i>Proceedings of the 2020 9th International Conference on Software and Computer Applications, ICSCA '20</i> , page 36–40, New York, NY, USA. Association for Computing Machinery.	
794		
795		
796		
797		
798		
799		
800	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	
801		
802		
803		
804		
805		
806		
807		
808	F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. <i>Journal of Machine Learning Research</i> , 12:2825–2830.	813
809		814
810		
811		
812		
	Bibek Poudel, Adam Cook, Sekou Traore, and Shehlah Ameli. 2024. Documint: Docstring generation for python using small language models . <i>CoRR</i> , abs/2405.10243.	815
		816
		817
		818
	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners .	819
		820
		821
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>J. Mach. Learn. Res.</i> , 21(1).	822
		823
		824
		825
		826
	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton-Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, and 6 others. 2023. Code llama: Open foundation models for code . <i>CoRR</i> , abs/2308.12950.	827
		828
		829
		830
		831
		832
		833
		834
	Mobashir Sadat and Cornelia Caragea. 2024. MSciNLI: A diverse benchmark for scientific natural language inference . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1610–1629, Mexico City, Mexico. Association for Computational Linguistics.	835
		836
		837
		838
		839
		840
		841
		842
	Alireza Salemi, Julian Killingback, and Hamed Zamani. 2025. ExPerT: Effective and explainable evaluation of personalized long-form text generation . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 17516–17532, Vienna, Austria. Association for Computational Linguistics.	843
		844
		845
		846
		847
		848
	Sergio Servantez, Joe Barrow, Kristian Hammond, and Rajiv Jain. 2024. Chain of logic: Rule-based reasoning with large language models . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 2721–2733, Bangkok, Thailand. Association for Computational Linguistics.	849
		850
		851
		852
		853
		854
	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank . In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.	855
		856
		857
		858
		859
		860
		861
		862
	Supriyono, Aji Prasetya Wibawa, Suyono, and Fachrul Kurniawan. 2024. Advancements in natural language processing: Implications, challenges, and future directions . <i>Telematics and Informatics Reports</i> , 16:100173.	863
		864
		865
		866
		867

868	Zhen Tan, Dawei Li, Song Wang, Alimohammad	Punta Cana, Dominican Republic. Association for	926
869	Beigi, Bohan Jiang, Amrita Bhattacharjee, Man-	Computational Linguistics.	927
870	sooreh Karami, Jundong Li, Lu Cheng, and Huan Liu.		
871	2024. Large language models for data annotation and	Alex Warstadt, Amanpreet Singh, and Samuel R. Bow-	928
872	synthesis: A survey . In <i>Proceedings of the 2024 Con-</i>	man. 2019. Neural network acceptability judgments .	929
873	<i>ference on Empirical Methods in Natural Language</i>	<i>Transactions of the Association for Computational</i>	930
874	<i>Processing</i> , pages 930–957, Miami, Florida, USA.	<i>Linguistics</i> , 7:625–641.	931
875	Association for Computational Linguistics.		
876	Erik F. Tjong Kim Sang and Fien De Meulder.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	932
877	2003. Introduction to the CoNLL-2003 shared task:	Chaumond, Clement Delangue, Anthony Moi, Pier-	933
878	Language-independent named entity recognition . In	ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,	934
879	<i>Proceedings of the Seventh Conference on Natural</i>	Joe Davison, Sam Shleifer, Patrick von Platen, Clara	935
880	<i>Language Learning at HLT-NAACL 2003</i> , pages 142–	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven	936
881	147.	Le Scao, Sylvain Gugger, and 3 others. 2020. Trans-	937
		formers: State-of-the-art natural language processing .	938
882	Weixi Tong and Tianyi Zhang, 2024. CodeJudge: Eval-	In <i>Proceedings of the 2020 Conference on Empirical</i>	939
883	uating code generation with large language models .	<i>Methods in Natural Language Processing: System</i>	940
884	In <i>Proceedings of the 2024 Conference on Empiri-</i>	<i>Demonstrations</i> , pages 38–45, Online. Association	941
885	<i>cal Methods in Natural Language Processing</i> , pages	for Computational Linguistics.	942
886	20032–20051, Miami, Florida, USA. Association for		
887	Computational Linguistics.	Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu.	943
888	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	2024. Benchmarking benchmark leakage in large	944
889	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	language models . <i>CoRR</i> , abs/2404.18824.	945
890	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal		
891	Azhar, Aurélien Rodriguez, Armand Joulin, Edouard	Weihao Xuan, Rui Yang, Heli Qi, Qingcheng Zeng,	946
892	Grave, and Guillaume Lample. 2023. Llama: Open	Yunze Xiao, Aosong Feng, Dairui Liu, Yun Xing, Jun-	947
893	and efficient foundation language models . <i>CoRR</i> ,	jue Wang, Fan Gao, Jinghui Lu, Yuang Jiang, Huitao	948
894	abs/2302.13971.	Li, Xin Li, Kunyu Yu, Ruihai Dong, Shangding	949
895	Muhammad Uzair-Ul-Haq, Davide Rigoni, and Alessan-	Gu, Yuekang Li, Xiaofei Xie, and 13 others. 2025.	950
896	dro Sperduti. 2025. Llms as data annotators:	MMLU-ProX: A multilingual benchmark for ad-	951
897	How close are we to human performance . <i>CoRR</i> ,	vanced large language model evaluation . In <i>Proceed-</i>	952
898	abs/2504.15022.	<i>ings of the 2025 Conference on Empirical Methods</i>	953
899	Oriol Vinyals, Charles Blundell, Timothy Lillicrap, ko-	<i>in Natural Language Processing</i> , pages 1513–1532,	954
900	ray kavukcuoglu, and Daan Wierstra. 2016. Match-	Suzhou, China. Association for Computational Lin-	955
901	ing networks for one shot learning . In <i>Advances in</i>	guistics.	956
902	<i>Neural Information Processing Systems</i> , volume 29.	Dayu Yang, Tianyang Liu, Daoan Zhang, Antoine	957
903	Curran Associates, Inc.	Simoulin, Xiaoyi Liu, Yuwei Cao, Zhaopu Teng, Xin	958
904	Alex Wang, Amanpreet Singh, Julian Michael, Felix	Qian, Grey Yang, Jiebo Luo, and Julian McAuley.	959
905	Hill, Omer Levy, and Samuel Bowman. 2018. GLUE:	2025a. Code to think, think to code: A survey on	960
906	A multi-task benchmark and analysis platform for nat-	code-enhanced reasoning and reasoning-driven code	961
907	ural language understanding . In <i>Proceedings of the</i>	intelligence in LLMs . In <i>Proceedings of the 2025</i>	962
908	<i>2018 EMNLP Workshop BlackboxNLP: Analyzing</i>	<i>Conference on Empirical Methods in Natural Lan-</i>	963
909	<i>and Interpreting Neural Networks for NLP</i> , pages	<i>guage Processing</i> , pages 2586–2616, Suzhou, China.	964
910	353–355, Brussels, Belgium. Association for Com-	Association for Computational Linguistics.	965
911	putational Linguistics.		
912	Wenhan Wang, Chenyuan Yang, Zhijie Wang, Yuheng	Haotong Yang, Yi Hu, Shijia Kang, Zhouchen Lin, and	966
913	Huang, Zhaoyang Chu, Da Song, Lingming Zhang,	Muhan Zhang. 2025b. Number cookbook: Number	967
914	An Ran Chen, and Lei Ma. 2025. TestEval: Bench-	understanding of language models and how to im-	968
915	marking large language models for test case gener-	prove it . In <i>The Thirteenth International Conference</i>	969
916	ation . In <i>Findings of the Association for Computa-</i>	<i>on Learning Representations, ICLR 2025, Singapore,</i>	970
917	<i>tional Linguistics: NAACL 2025</i> , pages 3547–3562,	<i>April 24–28, 2025</i> . OpenReview.net.	971
918	Albuquerque, New Mexico. Association for Compu-		
919	tational Linguistics.	Jian Yang, Jiayi Yang, Wei Zhang, Jin Ke, Yibo Miao,	972
920	Yue Wang, Weishi Wang, Shafiq Joty, and Steven C.H.	Lei Zhang, Liqun Yang, Zeyu Cui, Yichang Zhang,	973
921	Hoi. 2021. CodeT5: Identifier-aware unified pre-	Zhoujun Li, Binyuan Hui, and Junyang Lin. 2025c.	974
922	trained encoder-decoder models for code understand-	CodeArena: Evaluating and aligning CodeLLMs on	975
923	ing and generation . In <i>Proceedings of the 2021</i>	human preference . In <i>Proceedings of the 2025 Con-</i>	976
924	<i>Conference on Empirical Methods in Natural Lan-</i>	<i>ference on Empirical Methods in Natural Language</i>	977
925	<i>guage Processing</i> , pages 8696–8708, Online and	<i>Processing</i> , pages 9683–9694, Suzhou, China. Asso-	978
		ciation for Computational Linguistics.	979
		Lin Yang, Chen Yang, Shutao Gao, Weijing Wang,	980
		Bo Wang, Qihao Zhu, Xiao Chu, Jianyi Zhou,	981
		Guangtai Liang, Qianxiang Wang, and Junjie Chen.	982

983 2024. [On the evaluation of large language models](#)
984 [in unit test generation](#). In *Proceedings of the 39th*
985 *IEEE/ACM International Conference on Automated*
986 *Software Engineering, ASE '24*, page 1607–1619,
987 New York, NY, USA. Association for Computing
988 Machinery.

989 Li Zhang, Liam Dugan, Hainiu Xu, and Chris Callison-
990 burch. 2023a. [Exploring the curious case of code](#)
991 [prompts](#). In *Proceedings of the 1st Workshop on*
992 *Natural Language Reasoning and Structured Expla-*
993 *nations (NLRSE)*, pages 9–17, Toronto, Canada. As-
994 sociation for Computational Linguistics.

995 Linghao Zhang, Jingshu Zhao, Chong Wang, and Peng
996 Liang. 2024. [Using Large Language Models for](#)
997 [Commit Message Generation: A Preliminary Study](#).
998 In *2024 IEEE International Conference on Software*
999 *Analysis, Evolution and Reengineering (SANER)*,
1000 pages 126–130, Los Alamitos, CA, USA. IEEE Com-
1001 puter Society.

1002 Ruoyu Zhang, Yanzeng Li, Yongliang Ma, Ming Zhou,
1003 and Lei Zou. 2023b. [LLMaAA: Making large lan-](#)
1004 [guage models as active annotators](#). In *Findings of the*
1005 *Association for Computational Linguistics: EMNLP*
1006 *2023*, pages 13088–13103, Singapore. Association
1007 for Computational Linguistics.

1008 Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.
1009 Character-level convolutional networks for text clas-
1010 sification. In *Proceedings of the 29th International*
1011 *Conference on Neural Information Processing Sys-*
1012 *tems - Volume 1, NIPS'15*, page 649–657, Cambridge,
1013 MA, USA. MIT Press.

1014 Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen,
1015 Siqu Zuo, Andrea Hu, Christopher A. Choquette-
1016 Choo, Jingyue Shen, Joe Kelley, Kshitij Bansal,
1017 Luke Vilnis, Mateo Wirth, Paul Michel, Peter Choy,
1018 Pratik Joshi, Ravin Kumar, Sarmad Hashmi, Shub-
1019 ham Agrawal, Zhitao Gong, and 7 others. 2024.
1020 [Codegemma: Open code models based on gemma.](#)
1021 *CoRR*, abs/2406.11409.

1022 Xin Zhou, Martin Weyssow, Ratnadira Widyasari, Ting
1023 Zhang, Junda He, Yunbo Lyu, Jianming Chang, Beiqi
1024 Zhang, Dan Huang, and David Lo. 2025. [Lessleak-](#)
1025 [bench: A first investigation of data leakage in llms](#)
1026 [across 83 software engineering benchmarks](#). *CoRR*,
1027 abs/2502.06215.

A Implementation Details 1028

Model access and licenses. The links (given
as URL) of the evaluated models in this paper are
listed in Table 10. 1029
1030
1031

OpenAI models (such as GPT-3.5-Turbo) are
governed by a proprietary license⁶, which permits
usage via their APIs for both businesses and
developers. DeepSeek uses a custom license for
its models that allows free distribution, repro-
duction and usage of model copies under certain
conditions (refer to Section 4 and Attachment
A from [https://github.com/deepseek-ai/](https://github.com/deepseek-ai/DeepSeek-Coder/blob/main/LICENSE-CODE)
[DeepSeek-Coder/blob/main/LICENSE-CODE](https://github.com/deepseek-ai/DeepSeek-Coder/blob/main/LICENSE-CODE))
which we meet. Qwen models are generally avail-
able under the permissive Apache 2.0 License⁷.
Llama and Gemma models use their custom
licenses⁸. For this study, we strictly adhere to the
terms and conditions defined by each license: we
access proprietary models via their respective APIs
and locally reproduce copies of the open-weight
models solely for evaluation purposes, without
modification or redistribution. All uses remain
compliant with the terms of use defined in the
aforementioned licenses, limiting our scope to
research. 1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052

Model deployment. For inference, we follow
the appropriate computation framework for each
provider: GPT and DeepSeek models are accessed
through their official APIs, while all other mod-
els are run locally: NL-LLMs using vLLM (Kwon
et al., 2023) and Code-LLMs using the Hugging-
Face transformers (Wolf et al., 2020) library. 1053
1054
1055
1056
1057
1058
1059

vLLM is widely used for efficient inference in re-
cent NL-LLM studies (Salemi et al., 2025; Feng
et al., 2025), but it primarily supports instruction-
mode generation. In contrast, code completion
is natively supported by transformers, which is
also the officially recommended framework for *all*
locally evaluated Code-LLMs in this work, includ-
ing Qwen Coder⁹, codegemma¹⁰ and CodeLlama¹¹.
To ensure fair comparison across models, we fur-
ther fix the maximum output length to 16 tokens,
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069

⁶<https://openai.com/policies/terms-of-use/>

⁷<https://github.com/QwenLM/Qwen?tab=Apache-2.0-1-ov-file>

⁸<https://www.llama.com/llama2/license/> and
<https://ai.google.dev/gemma/terms>

⁹<https://github.com/QwenLM/Qwen3-Coder/blob/main/README.md>

¹⁰<https://huggingface.co/google/codegemma-7b>

¹¹<https://huggingface.co/meta-llama/CodeLlama-7b-hf>

which matches the default setting in vLLM¹² and comfortably covers the length of all target labels used in our tasks (see Table 1).

Computation infrastructure. To ensure full-precision inference, models with at least 30B parameters (e.g., the 30B and 32B Qwen models) are run on a single NVIDIA H200 GPU, while smaller models are run on a single NVIDIA A100 GPU.

B Detailed Experimental Results

Table 9 spotlights the zero-shot performance of five Code-LLMs across the English and Chinese splits from the XNLI dataset.

Table 12 and Table 13 present the performance and redundancy scores of few-shot experiments, respectively.

Model	English	Chinese	Drop by (%)
DeepSeek-V3.1 (Code)	0.756	0.743	-1.72
Qwen3-Coder-30B	0.842	0.760	-9.74
codegemma-7b	0.572	0.504	-11.89
GPT-3.5 Turbo (Code)	0.680	0.599	-11.91
CodeLlama-13b	0.576	0.492	-14.58

Table 9: Comparison of zero-shot accuracy between English and Chinese subsets of XNLI. The drop ratio is measured relative to the English score, i.e. $(score_{ZH} - score_{EN})/score_{EN}$.

C Effect of Components in Code Prompt

To analyze the contribution of individual components in our code-style prompt, we conduct an ablation study with four prompt variants on the HCC-Staging dataset. We choose HCC-Staging because, as discussed in the Limitations section, it is newly introduced and thus minimizes the risk of data leakage during LLM pre-training. The four prompt variants are defined as follows:

full: the original code prompt (see the Python prompt in Appendix E.10 for example);

-typing: the full prompt with type hints (e.g., `Literal["True", "False"]`) removed;

-docstring: the full prompt with the docstring removed;

-both: the full prompt with both type hints and the docstring removed. The resulting prompt contains only the function signature (i.e., `def HCC_staging(text: str):`) followed by an empty function body (`pass`).

¹²https://docs.vllm.ai/en/v0.8.4/api/inference_params.html

Model	URL
GPT-3.5-Turbo	https://platform.openai.com/docs/models/gpt-3.5-turbo
DeepSeek-V3.1	https://api.deepseek.com/beta (API until 2025/09/21) and https://huggingface.co/deepseek-ai/DeepSeek-V3.1
Qwen3-Coder-30B	https://huggingface.co/Qwen/Qwen3-Coder-30B-A3B-Instruct
Qwen3-30B-Instruct	https://huggingface.co/Qwen/Qwen3-30B-A3B-Instruct-2507
Qwen2.5-Coder-32B	https://huggingface.co/Qwen/Qwen2.5-Coder-14B
Qwen2.5-32B-Instruct	https://huggingface.co/Qwen/Qwen2.5-32B-Instruct
Qwen2.5-Coder-14B	https://huggingface.co/Qwen/Qwen2.5-Coder-14B
Qwen2.5-14B-Instruct	https://huggingface.co/Qwen/Qwen2.5-14B-Instruct
Qwen2.5-Coder-7B	https://huggingface.co/Qwen/Qwen2.5-Coder-7B
Qwen2.5-7B-Instruct	https://huggingface.co/Qwen/Qwen2.5-7B-Instruct
Qwen2.5-Coder-3B	https://huggingface.co/Qwen/Qwen2.5-Coder-3B
Qwen2.5-3B-Instruct	https://huggingface.co/Qwen/Qwen2.5-3B-Instruct
CodeLlama-13b	https://huggingface.co/meta-llama/CodeLlama-13b-hf
CodeLlama-13b-Python	https://huggingface.co/meta-llama/CodeLlama-13b-Python-hf
Llama-2-13b	https://huggingface.co/meta-llama/Llama-2-13b-hf
CodeLlama-7b	https://huggingface.co/meta-llama/CodeLlama-7b-hf
CodeLlama-7b-Python	https://huggingface.co/meta-llama/CodeLlama-7b-Python-hf
Llama-2-7b	https://huggingface.co/meta-llama/Llama-2-7b-hf
codegemma-7b	https://huggingface.co/google/codegemma-7b
gemma-7b	https://huggingface.co/google/gemma-7b
codegemma-2b	https://huggingface.co/google/codegemma-2b
gemma-2b	https://huggingface.co/google/gemma-2b

Table 10: The URLs of evaluated models.

Table 11 reports the zero-shot accuracy and redundancy of the four prompt variants. The evaluated models are the best performing ones in general across 10 tasks (see Section 4.3 for selection criteria). (1) **Accuracy:** removing either type hints or the docstring leads to a comparable and modest accuracy change (less than 0.03) from the full prompt. In contrast, removing both components (i.e., leaving only the function name and its argument) results in zero accuracy. This indicates a low risk of data leakage and confirms that the task cannot be solved from prior knowledge alone, but instead requires explicit rule specification provided by the prompt. (2) **Redundancy:** removing either component has no effect on redundancy for DeepSeek-V3.1, GPT-3.5 Turbo, and CodeGemma-7B, but increases redundancy for CodeLlama. Overall, the full prompt consistently achieves the best balance between accuracy and redundancy across models.

Model	full	-typing	-docstring	-both
<i>Accuracy</i>				
DeepSeek-V3.1 (Code)	83.9	80.0	83.9	0.0
GPT-3.5 Turbo (Code)	58.1	64.5	61.3	0.0
Qwen3-Coder-30B	67.7	64.5	58.1	0.0
codegemma-7b	67.7	70.9	64.5	0.0
CodeLlama-13b	51.6	61.2	54.8	0.0
Avg.	65.8	68.3	64.5	0.0
<i>Redundancy</i>				
DeepSeek-V3.1 (Code)	0.0	0.0	0.0	100.0
GPT-3.5 Turbo (Code)	0.0	0.0	0.0	100.0
Qwen3-Coder-30B	30.9	29.8	30.4	100.0
codegemma-7b	0.0	0.0	0.0	100.0
CodeLlama-13b	0.0	91.3	91.2	100.0
Avg.	6.2	24.2	24.3	100.0

Table 11: Comparison of the four prompt variants in zero-shot results (%) over the HCC-Staging dataset.

D Annotation Details

For the HCC-Staging dataset, a specialized physician from a German university hospital extracted tumor size information following the instruction: *Reports more than 12 months prior to the tumor board meeting are excluded. Only tumor size information (number + unit) is extracted, and multiple tumor sizes may be annotated for a single report.* The annotator was compensated according to the hospital’s tariff, which is above the local minimum wage. The annotation process required approximately 12 hours in total, distributed over several weeks, with no more than 6 working hours on any single day, fully complying with local labor regulations.

E Prompts

E.1 SST-2

Natural language prompt.

You are a data annotator for sentiment classification. Your task is to read the following text and classify the text into either positive or negative based on its sentiment.

Please answer directly with "positive" or "negative" only.

{few_shot_examples}

Here is the text to classify:

Input: *{text}*

Output:

Few-shot example in natural language:

Input: *{text}*

Output: *{label}*

Code prompt (Python).

```
def classify_sentiment(text: str) -> Literal["positive", "negative"]:
    """Classify the sentiment of the given text as either positive or negative.
    Args:
        - text (str): The text to classify.
    Returns:
        Literal["positive", "negative"]: The sentiment label of the text.
    """
    pass
```

{few_shot_examples}

Test case for inference

text = *{text}*

assert (classify_sentiment(text) == "

Few-shot example in Python:

Example test case

text = *{text}*

assert (classify_sentiment(text) == *{label}*)"

Code prompt (JavaScript).

```
const assert = require('assert');
function classifySentiment(text) {
    /*
```

Classify the sentiment of the given text as either positive or negative.

Args:

- text (string): The text to classify.

Returns:

"positive" | "negative": The sentiment label of the text.

*/

// TODO: Implement classification logic

}

{few_shot_examples}

Test case for inference

const text = *{text}*;

assert.strictEqual(classifySentiment(text), "

Few-shot example in JavaScript:

Example test case

const text = *{text}*;

assert.strictEqual(classifySentiment(text), *{label}*)"

Code prompt (C++).

```
#include <cassert>
```

```
#include <string>
```

```
std::string classify_sentiment(const std::string& text) {
```

```
    /*
```

Metric # Shots	TC			NLI			RC	NR		\bar{X}
	SST-2	AG News	CoLA	MRPC	XNLI	MSciNLI	SciERC	Iris	HCC-Staging	
	Acc 4	Acc 4	MCC 4	Acc 4	Acc 4	Macro-F1 4	Micro-F1 2	Acc 4	Acc 4	
GPT-3.5 Turbo (Code)	93.7±0.24	63.4 ±1.42	61.1±1.73	76.8±0.48	71.5±0.57	42.5 ±0.82	60.8 ±2.80	96.0±10.97	58.1±0.54	69.3
GPT-3.5 Turbo (NL)	94.6 ±0.19	63.4 ±0.20	61.9 ±1.53	78.2 ±0.52	75.5 ±1.05	28.1±1.08	55.3±3.05	96.7 ±4.02	74.2 ±0.54	69.8
DeepSeek-V3.1 (Code)	94.7±0.25	63.1 ±0.51	71.0 ±1.30	78.0±0.44	85.1±0.35	47.4±1.93	33.7±0.51	96.0 ±3.04	87.1±0.00	72.9
DeepSeek-V3.1 (NL)	95.8 ±0.69	63.1 ±0.74	70.5±0.28	80.2 ±0.43	87.4 ±0.40	63.4 ±1.80	66.7 ±4.02	96.0 ±1.52	90.3 ±0.00	79.3 †
Qwen3-Coder-30B	95.3 ±0.19	65.1 ±0.35	60.4±1.43	79.0 ±0.50	86.8 ±0.25	57.4±1.67	53.9±3.20	95.7 ±1.52	64.5±0.68	73.1
Qwen3-30B-Instruct	88.3±5.30	61.2±1.91	64.1 ±0.74	42.2±11.91	85.5±1.49	60.8 ±3.42	67.2 ±2.11	90.1±6.08	77.4 ±3.69	70.8
Qwen2.5-Coder-32B	95.8 ±0.28	63.9 ±1.87	59.3±1.17	72.8±0.90	87.2±0.22	53.5 ±0.61	42.0±4.59	97.8 ±3.04	83.9 ±2.08	72.9
Qwen2.5-32B-Instruct	94.4±0.14	63.9 ±0.32	60.3 ±1.41	78.6 ±0.56	89.6 ±1.04	47.0±4.34	68.0 ±2.06	95.7±10.54	64.5±0.00	73.6
Qwen2.5-Coder-14B	94.6±0.24	59.8±1.32	53.3±2.38	71.7±1.19	85.5±0.79	50.7 ±1.53	52.6±0.35	92.8 ±4.02	77.4 ±1.77	70.9
Qwen2.5-14B-Instruct	95.4 ±0.61	63.6 ±0.18	59.5 ±0.38	75.7 ±1.43	85.9 ±0.69	43.4±4.02	63.5 ±6.11	91.3±1.52	71.0±3.13	72.1
Qwen2.5-Coder-7B	93.5±0.52	52.0±5.17	50.2±3.52	71.4 ±0.76	82.7±1.34	35.6±5.94	26.4±1.38	88.4±5.27	74.2 ±3.13	63.8
Qwen2.5-7B-Instruct	94.0 ±0.24	55.2 ±1.61	58.5 ±1.76	51.4±4.52	83.5 ±0.71	41.8 ±1.35	45.5 ±5.10	95.7 ±5.48	61.3±3.04	65.2
Qwen2.5-Coder-3B	92.1 ±0.27	47.8±3.13	45.8±1.91	71.2 ±0.41	78.4±0.94	39.0±1.56	30.1±1.65	79.7±12.07	77.4 ±1.81	62.4
Qwen2.5-3B-Instruct	89.3±2.51	57.1 ±2.40	51.3 ±1.14	39.9±0.36	81.9 ±0.92	46.3 ±0.60	48.4 ±2.38	89.1 ±10.64	58.1±1.71	62.4
CodeLlama-13b	91.9 ±1.19	49.0±1.38	35.8±2.48	72.1±0.49	60.0 ±0.48	42.2±1.52	16.2±2.36	79.7±8.47	51.6 ±2.13	55.4
CodeLlama-13b-Python	88.9±1.59	52.7 ±9.30	40.7 ±4.59	75.1 ±3.34	58.5±0.47	43.0 ±2.01	16.8 ±1.30	80.4 ±12.07	48.4±4.16	56.1
Llama-2-13b	0.0±0.00	0.9±1.10	0.0±0.00	68.7 ±1.32	41.1±1.38	21.4±4.36	5.7±0.58	0.0±0.00	0.0±0.00	15.3
CodeLlama-7b	93.7 ±1.70	52.8±1.88	24.4 ±0.87	71.0±0.37	55.8 ±0.77	30.4±0.80	17.8±1.21	83.3±11.88	67.7 ±2.08	55.2
CodeLlama-7b-Python	91.9±1.52	53.3 ±1.01	9.1±1.66	75.1 ±0.18	54.0±0.44	36.0 ±3.44	19.9 ±1.56	90.6 ±9.50	64.5±1.49	54.9
Llama-2-7b	0.0±0.43	1.9±9.40	1.0±3.17	23.9 ±2.36	0.7±0.29	11.3±0.61	18.0±3.50	0.0±1.52	0.0±0.00	6.3
codegemma-7b	94.0 ±0.43	62.7 ±2.16	40.3 ±7.50	54.8 ±5.60	66.1 ±1.85	48.7 ±4.08	20.8 ±2.60	91.3 ±2.63	80.6 ±1.02	62.1
gemma-7b	1.4±1.55	1.8±0.60	0.0±0.96	39.2±16.12	56.7±2.01	23.7±6.45	19.0±0.50	0.0±3.04	19.4±0.00	17.9
codegemma-2b	71.4 ±4.58	25.8 ±0.28	3.4 ±1.88	66.5 ±0.00	39.1 ±1.19	10.0±0.00	9.3±1.17	64.5 ±13.69	71.0 ±3.47	40.1
gemma-2b	5.7±8.58	1.0±0.64	2.4±2.80	63.4±1.29	29.8±1.02	17.8 ±2.90	12.3 ±3.28	33.3±5.20	48.4±0.00	23.8
Avg. Code-LLM	91.7	54.7	42.7	72.0	70.1	41.3	30.8	87.4	69.7	62.2
Avg. NL-LLM	59.9	39.4	39.0	58.3	65.2	36.8	42.7	55.3	51.3	49.8
Avg. NL-LLM (w/o llama, gemma)	93.1	61.1	60.9	63.7	84.2	47.3	59.2	82.1	71.0	69.2

Table 12: Few-shot performance (%) of NL-LLMs using natural language prompts and Code-LLMs using code-based prompts. Both mean and standard deviation across three runs are reported. The higher score within each comparison is highlighted in bold. The highest score of each task is denoted with †. Acc: Accuracy. MCC: Matthews correlation coefficient. \bar{X} denotes the arithmetic mean across 10 tasks.

# Shots	TC			NLI			RC	NR		\bar{R}
	SST-2	AG News	CoLA	MRPC	XNLI	MSciNLI	SciERC	Iris	HCC-Staging	
	4	4	4	4	4	4	2	4	4	
GPT-3.5 Turbo (Code)	0.0 ±0.00	25.2±0.04	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.02	0.3 ±0.11	0.0 ±0.00	0.0 ±0.00	2.8
GPT-3.5 Turbo (NL)	0.0 ±0.00	25.1 ±0.14	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	10.3±0.12	0.0 ±0.00	0.0 ±0.00	3.9
DeepSeek-V3.1 (Code)	0.0 ±0.00	25.3±0.09	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	94.5±0.00	0.0 ±0.00	0.0 ±0.00	13.3
DeepSeek-V3.1 (NL)	0.0 ±0.00	25.4 ±0.03	2.6±0.00	0.0 ±0.00	10.5±0.00	6.8±0.00	0.0 ±0.00	2.4±0.00	0.0 ±0.00	5.3
Qwen3-Coder-30B	19.3 ±0.55	38.6 ±1.07	15.9 ±0.15	14.3 ±1.86	17.6 ±0.15	18.5 ±0.14	14.7 ±0.59	19.9 ±0.33	29.5 ±0.01	20.9
Qwen3-30B-Instruct	88.8±1.68	92.6±0.79	86.1±0.31	79.6±1.13	84.0±0.28	88.1±0.31	54.1±5.86	98.6±0.27	92.1±0.52	84.9
Qwen2.5-Coder-32B	0.0 ±0.00	25.7 ±0.18	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	2.9
Qwen2.5-32B-Instruct	88.8±0.25	91.9±0.16	81.1±1.73	83.8±0.94	83.8±0.05	88.4±0.61	86.0±0.60	68.8±0.15	91.5±3.18	84.9
Qwen2.5-Coder-14B	0.0 ±0.00	25.7 ±1.10	0.0 ±0.00	0.0 ±0.00	0.0 ±0.02	2.0 ±7.59	1.0 ±1.26	32.9 ±0.00	0.0 ±0.00	6.8
Qwen2.5-14B-Instruct	86.2±11.22	92.1±0.16	85.4±0.24	82.5±0.61	85.9±0.44	88.7±0.23	88.6±0.47	81.5±0.41	91.3±0.27	86.9
Qwen2.5-Coder-7B	0.0 ±0.00	20.4 ±0.96	0.0 ±0.00	0.0 ±0.00	68.8 ±9.62	7.2 ±13.26	0.8 ±0.42	0.0 ±0.00	0.0 ±0.00	10.8
Qwen2.5-7B-Instruct	88.4±0.07	91.7±0.12	85.2±0.23	68.4±4.98	85.3±0.21	88.3±0.37	89.4±0.56	88.2±0.59	89.9±0.26	86.1
Qwen2.5-Coder-3B	4.5 ±5.50	15.9 ±3.82	0.0 ±0.30	0.0 ±0.05	80.7 ±9.88	6.1 ±9.98	6.5 ±3.83	23.2 ±1.41	3.0 ±2.11	15.5
Qwen2.5-3B-Instruct	59.8±6.68	92.3±0.70	81.7±1.07	68.4±3.82	84.8±0.27	86.8±0.97	86.9±0.21	82.8±0.16	93.1±0.32	81.8
CodeLlama-13b	84.6±0.07	89.8±0.14	80.0±0.52	79.3 ±0.24	78.3±0.30	80.1 ±0.27	80.5±0.51	83.0±0.18	91.8±0.14	83.0
CodeLlama-13b-Python	84.8±0.09	88.9 ±1.06	79.8 ±0.30	79.3 ±0.51	76.8 ±0.11	82.1±0.22	74.6 ±0.67	82.4±0.24	91.7 ±0.53	82.3
Llama-2-13b	0.0 ±0.00	99.8±0.19	100.0±0.00	82.8±3.27	82.1±1.87	96.5±1.27	95.0±1.04	0.0 ±0.00	100.0±0.00	72.9
CodeLlama-7b	84.5±0.13	89.2±0.61	77.8±0.73	79.9±0.07	77.3±0.06	79.8 ±0.38	75.7 ±1.43	82.7±0.23	91.4±0.14	82.0
CodeLlama-7b-Python	60.3 ±10.49	87.7 ±3.30	75.2 ±2.37	79.0±0.32	76.7 ±0.23	80.4±0.27	77.3±1.55	15.3±7.71	52.1 ±1.73	67.1
Llama-2-7b	100.0±0.12	99.7±2.97	100.0±0.62	73.3 ±0.49	99.5±0.19	99.6±0.18	92.7±1.47	0.0 ±0.13	100.0±0.00	85.0
codegemma-7b	0.0 ±0.00	24.3 ±0.18	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	2.7
gemma-7b	99.6±0.42	99.2±0.26	95.9±0.58	89.2±4.70	85.1±0.50	93.6±3.01	80.8±1.64	100.0±19.86	36.5±0.00	86.7
codegemma-2b	4.4 ±13.71	2.7 ±0.51	82.6 ±0.23	40.7 ±17.86	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	0.0 ±0.00	6.0±0.00	15.2
gemma-2b	98.9±1.12	99.1±0.56	87.4±2.35	83.2±2.18	82.2±0.75	95.7±2.42	83.3±2.66	83.3±18.34	2.6 ±0.00	79.5
Avg. Code-LLM	26.3	43.0	31.6	28.7	36.6	27.4	32.8	26.1	28.1	31.2
Avg. NL-LLM	64.6	82.6	73.2	64.7	71.2	75.7	69.7	55.1	63.4	68.9
Avg. NL-LLM (w/o llama, gemma)	58.9	73.0	60.3	54.7	62.0	63.9	59.3	60.3	65.4	62.0

Table 13: Redundancy (%) (defined in Section 3.3; lower values indicate better generation precision) of few-shot experiments. Mean and standard deviation of redundancy across three runs are reported. The better (i.e. lower) redundancy in each comparison is highlighted in bold.

1208	Classify the sentiment of the given text as either	pass	1259
1209	positive or negative.		1260
1210	Args:	<i>{few_shot_examples}</i>	1261
1211	- text (std::string): The text to classify.	# Test case for inference	1262
1212	Returns:	text = "{text}"	1263
1213	"positive" or "negative": The sentiment label	assert (classify_topic(text) == "	1264
1214	of the text.	Few-shot example in Python:	1265
1215	*/	# Example test case	1266
1216	// TODO: Implement classification logic	text = "{text}"	1267
1217	}	assert (classify_topic(text) == "{label}")	1268
1218	<i>{few_shot_examples}</i>		1269
1219	# Test case for inference		
1220	std::string text = "{text}";		
1221	assert(classify_sentiment(text) == "		
1222	Few-shot example in C++:		
1223	# Example test case		
1224	std::string text = "{text}";		
1225	assert(classify_sentiment(text) == "{label}")		
1226			
1227	E.2 AG News		
1228	Natural language prompt.	E.3 CoLA	1270
1229	You are a data annotator for topic classification of	Natural language prompt.	1271
1230	news. Your task is to read the following news (Title	You are a linguist to annotate grammatical accep-	1272
1231	+ Description) and classify the news into one of the	tance. Your task is to read the following text and	1273
1232	topics: World, Sport, Business, or Sci/Tech.	tell if it is grammatically acceptable or unaccept-	1274
1233	Please answer directly with "World", "Sport",	able in standard English.	1275
1234	"Business", or "Sci/Tech" only.	Please answer with "acceptable" or "unacceptable"	1276
1235	<i>{few_shot_examples}</i>	only.	1277
1236	Here is the text to classify:	<i>{few_shot_examples}</i>	1278
1237	Input: {text}	Here is the text to classify:	1279
1238	Output:	Input: {text}	1280
1239	Few-shot example in natural language:	Output:	1281
1240	Input: {text}	Few-shot example in natural language:	1282
1241	Output: {label}	Input: {text}	1283
1242		Output: {label}	1284
1243	Code prompt (Python).		1285
1244	def classify_topic(text: str) -> Literal["World",	Code prompt (Python).	1286
1245	"Sport", "Business", "Sci/Tech"]:	def classify_grammatical_acceptance(text: str) ->	1287
1246	"""Classify the given text (title + description of	Literal["acceptable", "unacceptable"]:	1288
1247	news) into one of the following topics:	"""Classify the given text as either grammati-	1289
1248	- World	cally "acceptable" or "unacceptable".	1290
1249	- Sport	n	1291
1250	- Business	n	1292
1251	- Sci/Tech	Args:	1293
1252	Args:	- text (str): The text to classify.	1294
1253	- text (str): The text of news, which consists	Returns:	1295
1254	of title and description.	Literal["acceptable", "unacceptable"]: The	1296
1255	Returns:	grammatical acceptance of the text.	1297
1256	Literal["World", "Sport", "Business",	"""	1298
1257	"Sci/Tech"]: The topic label of the text.	pass	1299
1258	"""		1300
		<i>{few_shot_examples}</i>	1301
		# Test case for inference	1302
		text = "{text}"	1303
		assert (classify_grammatical_acceptance(text) ==	1304
		"	1305
		Few-shot example in Python:	1306
		# Example test case	1307
		text = "{text}"	1308
		assert (classify_grammatical_acceptance(text) ==	1309

1310	<code>"{label}")</code>	<code>== "{label}")</code>	1361
1311			1362
1312	E.4 MRPC	Code (JavaScript) prompt.	1363
1313	Natural language prompt.	<code>const assert = require('assert');</code>	1364
1314	You are a linguist to annotate sentence pairs for	<code>function detectParaphrase(sentence1, sentence2) {</code>	1365
1315	semantic equivalence. Your task is to read two sen-	<code>/*</code>	1366
1316	tences and determine whether they are paraphrases	Detect if two sentences are paraphrases of each	1367
1317	of each other or not.	other.	1368
1318	Paraphrase (equivalent) means the two sentences	Paraphrase means the two sentences express the	1369
1319	express the same meaning. Not paraphrase	same meaning.	1370
1320	(not_equivalent) means they do not express the	Args:	1371
1321	same meaning.	- sentence1 (string): The first sentence.	1372
1322	Please answer with one word only: "equivalent" or	- sentence2 (string): The second sentence.	1373
1323	"not_equivalent".	Returns:	1374
1324	<code>{few_shot_examples}</code>	"equivalent" or "not_equivalent": "equivalent"	1375
1325	Sentence 1: <code>{sentence1}</code>	if the sentences are paraphrases, "not_equivalent"	1376
1326	Sentence 2: <code>{sentence2}</code>	otherwise.	1377
1327	Relation:	*/	1378
1328	Few-shot example in natural language:	// TODO: Implement paraphrase detection logic	1379
1329	Sentence 1: <code>{sentence1}</code>	<code>}</code>	1380
1330	Sentence 2: <code>{sentence2}</code>	<code>{few_shot_examples}</code>	1381
1331	Relation: <code>{label}</code>	// Test case for inference	1382
1332		<code>const sentence1 = "{sentence1}";</code>	1383
1333	Code prompt (Python).	<code>const sentence2 = "{sentence2}";</code>	1384
1334	<code>def detect_paraphrase(sentence1: str, sentence2:</code>	<code>assert.strictEqual(detectParaphrase(sentence1, sen-</code>	1385
1335	<code>str) -> Literal["equivalent", "not_equivalent"]:</code>	<code>tence2), "</code>	1386
1336	"""Detect if two sentences are paraphrases of	Few-shot example in JavaScript:	1387
1337	each other.	# Example test case	1388
1338	Paraphrase means the two sentences express the	<code>const sentence1 = "{sentence1}";</code>	1389
1339	same meaning.	<code>const sentence2 = "{sentence2}";</code>	1390
1340	Args:	<code>assert.strictEqual(detectParaphrase(sentence1,</code>	1391
1341	- sentence1 (str): The first sentence.	<code>sentence2), "{label}");</code>	1392
1342	- sentence2 (str): The second sentence.		1393
1343	Returns:	Code (C++) prompt.	1394
1344	Literal["equivalent", "not_equivalent"]:	<code>#include <cassert></code>	1395
1345	"equivalent" if the sentences are paraphrases,	<code>#include <string></code>	1396
1346	"not_equivalent" otherwise.	<code>std::string detect_paraphrase(const std::string&</code>	1397
1347	"""	<code>sentence1, const std::string& sentence2) {</code>	1398
1348	pass	/*	1399
1349		Detect if two sentences are paraphrases of each	1400
1350	<code>{few_shot_examples}</code>	other.	1401
1351	# Test case for inference	Paraphrase means the two sentences express the	1402
1352	<code>sentence1 = "{sentence1}"</code>	same meaning.	1403
1353	<code>sentence2 = "{sentence2}"</code>	Args:	1404
1354	<code>assert (detect_paraphrase(sentence1, sentence2)</code>	- sentence1 (std::string): The first sentence.	1405
1355	<code>== "</code>	- sentence2 (std::string): The second sentence.	1406
1356	Few-shot example in Python:	Returns:	1407
1357	# Example test case	"equivalent" or "not_equivalent": "equivalent"	1408
1358	<code>sentence1 = "{sentence1}"</code>	if the sentences are paraphrases, "not_equivalent"	1409
1359	<code>sentence2 = "{sentence2}"</code>	otherwise.	1410
1360	<code>assert (detect_paraphrase(sentence1, sentence2)</code>	*/	1411
		// TODO: Implement paraphrase detection logic	1412

1515 sentences extracted from publications in the field
1516 of computer science. Your task is to read sentence1
1517 and sentence2, and then determine whether the
1518 relation between them is entailment, contrasting,
1519 reasoning or neutral.
1520 Entailment means sentence2 can be appended to
1521 sentence1 using words like "Specifically", "Pre-
1522 cisely", "In particular", "Particularly", "That is",
1523 "In other words".
1524 Contrasting means sentence2 can be appended to
1525 sentence1 using words like "However", "On the
1526 other hand", "In contrast", "On the contrary".
1527 Reasoning means sentence2 can be appended to
1528 sentence1 using words like "Therefore", "Thus",
1529 "Consequently", "As a result", "As a consequence",
1530 "From here we can infer".
1531 Neutral means none of the above.
1532 Please answer with one word only: "entailment",
1533 "contrasting", "reasoning" or "neutral".
1534 *{few_shot_examples}*
1535 sentence1: *{sentence1}*
1536 sentence2: *{sentence2}*
1537 relation:
1538 **Few-shot example in natural language:**
1539 sentence1: *{sentence1}*
1540 sentence2: *{sentence2}*
1541 relation: *{label}*
1542
1543 **Code (Python) prompt.**
1544 def natural_language_inference(sentence1: str,
1545 sentence2: str) -> Literal["entailment", "contrast-
1546 ing", "reasoning", "neutral"]:
1547 """Classify the relation between sentence1 and
1548 sentence2 from publications
1549 in the field of computer science into one of 4
1550 categories:
1551 "entailment" means sentence2 can be appended
1552 to sentence1 using words like
1553 'Specifically', 'Precisely', 'In particular',
1554 'Particularly', 'That is', 'In other words'.
1555 "contrasting" means sentence2 can be appended
1556 to sentence1 using words like
1557 'However', 'On the other hand', 'In contrast',
1558 'On the contrary'.
1559 "reasoning" means sentence2 can be appended
1560 to sentence1 using words like
1561 'Therefore', 'Thus', 'Consequently', 'As a
1562 result', 'As a consequence', 'From here we can
1563 infer'.
1564 "neutral" means none of the above.
1565 Args:
1566 - sentence1 (str): The first sentence.

- sentence2 (str): The second sentence. 1567
Returns: 1568
Literal["entailment", "contrasting", "reason- 1569
ing", "neutral"]: The relation label. 1570
""" 1571
pass 1572

{few_shot_examples} 1573
Test case for inference 1574
sentence1 = "*{sentence1}*" 1575
sentence2 = "*{sentence2}*" 1576
assert (natural_language_inference(sentence1, 1577
sentence2) == " 1578
Few-shot example in Python: 1579
Example test case 1580
sentence1 = "*{sentence1}*" 1581
sentence2 = "*{sentence2}*" 1582
assert (natural_language_inference(sentence1, 1583
sentence2) == "*{label}*") 1584
1585
1586

E.7 SemEval 1587

Natural language prompt. 1588

You are a data annotator for relation classification 1589
of two entities in a text. Your task is to read the 1590
following text in which two entities are marked by 1591
<e1> and </e1> for the subject entity, and <e2> 1592
and </e2> for the object entity, and classify the 1593
relation between the two entities into one of the 1594
following 10 labels: 1595
"Cause-Effect", "Component-Whole", "Content- 1596
Container", "Entity-Destination", "Entity-Origin", 1597
"Instrument-Agency", "Member-Collection", 1598
"Message-Topic", "Product-Producer", "Other", 1599
where "Other" means there is no relation between 1600
the two entities. 1601

{few_shot_examples} 1602
Here is the text for relation classification: 1603
Input: *{text}* 1604
Output: The relation between *{subj}* and *{obj}* in 1605
the input is 1606
Few-shot example in natural language: 1607
Input: *{text}* 1608
Output: The relation between *{subj}* and *{obj}* in 1609
the input is *{label}*. 1610
1611

Code (Python) prompt. 1612

def classify_relation(text: str, subj: str, obj: str) 1613
-> Literal["Cause-Effect", "Component-Whole", 1614
"Content-Container", "Entity-Destination", 1615
"Entity-Origin", "Instrument-Agency", "Member- 1616
Collection", "Message-Topic", "Product-Producer", 1617

```

1618 "Other"]:
1619     """"Classify the relation of the subject entity and
1620 object entity from the given text
1621 into one of the following 10 labels:
1622 "Cause-Effect", "Component-Whole", "Content-
1623 Container", "Entity-Destination",
1624 "Entity-Origin", "Instrument-Agency",
1625 "Member-Collection", "Message-Topic",
1626 "Product-Producer", "Other", where "Other"
1627 means there is no relation between the two entities.
1628 Args:
1629     - text (str): A text with a subject entity marked
1630 by <e1> and </e1>, and an object entity marked by
1631 <e2> and </e2>.
1632     - subj (str): The subject entity in the text.
1633     - obj (str): The object entity in the text.
1634 Returns:
1635     Literal["Cause-Effect", "Component-Whole",
1636 "Content-Container", "Entity-Destination",
1637 "Entity-Origin", "Instrument-Agency", "Member-
1638 Collection", "Message-Topic", "Product-Producer",
1639 "Other"]: The relation label between the subject
1640 entity and object entity.
1641     """"
1642     pass
1643
1644 {few_shot_examples}
1645 # Test case for inference
1646 text = "{text}"
1647 subj, obj = "{subj}", "{obj}"
1648 assert (classify_relation(text, subj, obj) == "
1649     Few-shot example in Python:
1650 # Example test case
1651 text = "{text}"
1652 subj, obj = "{subj}", "{obj}"
1653 assert (classify_relation(text, subj, obj) == "{la-
1654 bel}")
1655
1656 Code (JavaScript) prompt.
1657 const assert = require('assert');
1658 function classifyRelation(text, subj, obj) {
1659     /*
1660     Classify the relation of the subject entity and
1661 object entity from the given text
1662 into one of the following 10 labels:
1663 "Cause-Effect", "Component-Whole", "Content-
1664 Container", "Entity-Destination",
1665 "Entity-Origin", "Instrument-Agency",
1666 "Member-Collection", "Message-Topic",
1667 "Product-Producer", "Other", where "Other"
1668 means there is no relation between the two entities.
1669 Args:

```

```

- text (string): A text with a subject entity
1670 marked by <e1> and </e1>, and an object entity
1671 marked by <e2> and </e2>.
1672
- subj (string): The subject entity in the text.
1673
- obj (string): The object entity in the text.
1674
Returns:
1675
"Cause-Effect" | "Component-Whole" |
1676 "Content-Container" | "Entity-Destination" |
1677 "Entity-Origin" | "Instrument-Agency" |
1678 "Member-Collection" | "Message-Topic" |
1679 "Product-Producer" | "Other": The relation
1680 label between the subject entity and object entity.
1681
*/
1682 // TODO: Implement classification logic
1683
}
1684
{few_shot_examples}
1685
## Test case for inference
1686
const text = "{text}";
1687
const subj = "{subj}", obj = "{obj}";
1688
assert.strictEqual(classifyRelation(text, subj, obj),
1689 "
1690
Few-shot example in JavaScript:
1691
# Example test case
1692
const text = "{text}";
1693
const subj = "{subj}", obj = "{obj}";
1694
assert.strictEqual(classifyRelation(text, subj, obj),
1695 "{label}");
1696
Code (C++) prompt.
1697
#include <cassert>
1698
#include <string>
1699
1700
std::string classify_relation(const std::string& text,
1701 const std::string& subj, const std::string& obj) {
1702
/*
1703
Classify the relation of the subject entity and
1704 object entity from the given text
1705 into one of the following 10 labels:
1706 "Cause-Effect", "Component-Whole", "Content-
1707 Container", "Entity-Destination",
1708 "Entity-Origin", "Instrument-Agency",
1709 "Member-Collection", "Message-Topic",
1710 "Product-Producer", "Other", where "Other"
1711 means there is no relation between the two entities.
1712 Args:
1713     - text (std::string): A text with a subject entity
1714 marked by <e1> and </e1>, and an object entity
1715 marked by <e2> and </e2>.
1716
1717     - subj (std::string): The subject entity in the
1718 text.
1719
1720     - obj (std::string): The object entity in the
1721 text.

```

```

1722     Returns:
1723         "Cause-Effect" or "Component-Whole" or
1724 "Content-Container" or "Entity-Destination" or
1725     "Entity-Origin" or "Instrument-Agency" or
1726 "Member-Collection" or "Message-Topic" or
1727     "Product-Producer" or "Other": The relation
1728 label between the subject entity and object entity.
1729     */
1730     // TODO: Implement classification logic
1731 }
1732 {few_shot_examples}
1733 ## Test case for inference
1734 std::string text = "{text}";
1735 std::string subj = "{subj}", obj = "{obj}";
1736 assert(classify_relation(text, subj, obj) == "
1737     Few-shot example in C++:
1738 # Example test case
1739 std::string text = "{text}";
1740 std::string subj = "{subj}", obj = "{obj}";
1741 assert(classify_relation(text, subj, obj) == "{la-
1742 bel}");
1743
1744 Code (Ruby) prompt.
1745 def classify_relation(text, subj, obj)
1746     # Classify the relation of the subject entity and
1747 object entity from the given text
1748     # into one of the following 10 labels:
1749     # "Cause-Effect", "Component-Whole",
1750 "Content-Container", "Entity-Destination",
1751     # "Entity-Origin", "Instrument-Agency",
1752 "Member-Collection", "Message-Topic",
1753     # "Product-Producer", "Other", where "Other"
1754 means there is no relation between the two entities.
1755     # Args:
1756     # - text (String): A text with a subject entity
1757 marked by <e1> and </e1>, and an object entity
1758 marked by <e2> and </e2>.
1759     # - subj (String): The subject entity in the
1760 text.
1761     # - obj (String): The object entity in the text.
1762     # Returns:
1763     # "Cause-Effect", "Component-Whole",
1764 "Content-Container", "Entity-Destination",
1765     # "Entity-Origin", "Instrument-Agency",
1766 "Member-Collection", "Message-Topic",
1767     # "Product-Producer", "Other": The relation
1768 label between the subject entity and object entity.
1769     # TODO: Implement classification logic
1770 end
1771 {few_shot_examples}
1772 # Test case for inference
1773 text = "{text}"

```

```

subj, obj = "{subj}", "{obj}"
1774 raise "Assertion failed" unless (clas-
1775 sify_relation(text, subj, obj) == "
1776     Few-shot example in Ruby:
1777 # Example test case
1778 text = "{text}"
1779 subj, obj = "{subj}", "{obj}"
1780 raise "Assertion failed" unless (clas-
1781 sify_relation(text, subj, obj) == "{label}")
1782
1783

```

E.8 SciERC

Natural language prompt.

You are a data annotator for relation classification of two entities in a text. Your task is to read the following text in which the subject entity is surrounded by '[' and ']', and object entity is surrounded by '<' and '>', and then classify the relation between the two entities into one of the following 7 labels: "compare", "conjunction", "evaluate-for", "feature-of", "hyponym-of", "part-of", "used-for".

{few_shot_examples}

Here is the text for relation classification:

Input: {text}

Output: The relation between {subj} and {obj} in the input is

Few-shot example in natural language:

Input: {text}

Output: The relation between {subj} and {obj} in the input is {label}

Code (Python) prompt.

```

def classify_relation(text: str, subj: str, obj: str) ->
1806 Literal["compare", "conjunction", "evaluate-for",
1807 "feature-of", "hyponym-of", "part-of", "used-for"]:
1808     """Classify the relation of the subject entity and
1809 object entity from the given text
1810     into one of the following 7 labels: "compare",
1811 "conjunction", "evaluate-for", "feature-of",
1812 "hyponym-of", "part-of", "used-for".
1813     Args:
1814     - text (str): A text with a subject entity
1815 surrounded by '[' and ']', and an object entity
1816 surrounded by '<' and '>'.
1817     - subj (str): The subject entity in the text.
1818     - obj (str): The object entity in the text.
1819     Returns:
1820     Literal["compare", "conjunction", "evaluate-
1821 for", "feature-of", "hyponym-of", "part-of",
1822 "used-for"]: The relation label between the subject
1823 entity and object entity.
1824

```

1825	"""	If Petal Width \geq 1.8 cm, then the species is	1876
1826	pass	"virginica".	1877
1827		Args:	1878
1828	{few_shot_examples}	- text (str): A string containing the features of	1879
1829	# Test case for inference	the iris flower in the format: "sepal length: X cm,	1880
1830	text = "{text}"	sepal width: Y cm, petal length: Z cm, petal width:	1881
1831	subj, obj = "{subj}", "{obj}"	W cm".	1882
1832	assert (classify_relation(text, subj, obj) == "	Returns:	1883
1833	Few-shot example in Python:	Literal["setosa", "versicolor", "virginica"]:	1884
1834	# Example test case	The species label of the iris flower.	1885
1835	text = "{text}"	"""	1886
1836	subj, obj = "{subj}", "{obj}"	pass	1887
1837	assert (classify_relation(text, subj, obj) == "{la-		1888
1838	bel}")		1889
1839	E.9 Iris		1890
1840	Natural language prompt.		1891
1841	You are a data annotator for species classification	text = "sepal length: {sepal_length} cm,	1892
1842	of iris flowers. Your task is to read the following	sepal width: {sepal_width} cm, petal length:	1893
1843	rule for classification and the numerical features of	{petal_length} cm, petal width: {petal_width} cm"	1894
1844	the iris flower, and classify the species into either	assert (classify_iris_species(text) == "	1895
1845	setosa, versicolor, or virginica based on the rule.	Few-shot example in Python:	1896
1846	Please answer directly with only one of: "setosa",	# Example test case	1897
1847	"versicolor", or "virginica" as output.	text = "sepal length: {sepal_length} cm,	1898
1848	Rule:	sepal width: {sepal_width} cm, petal length:	1899
1849	If Petal Width $<$ 0.8 cm, then the species is	{petal_length} cm, petal width: {petal_width} cm"	1900
1850	"setosa".	assert (classify_iris_species(text) == "{label}")	1901
1851	If Petal Width \geq 0.8 cm and $<$ 1.8 cm, then the	E.10 HCC-Staging	1902
1852	species is "versicolor".	Natural language prompt.	1903
1853	If Petal Width \geq 1.8 cm, then the species is	You are a clinical annotator to determine whether a	1904
1854	"virginica".	patient with hepatocellular carcinoma (HCC) meets	1905
1855	{few_shot_examples}	the Milan criteria based on reported tumor size(s).	1906
1856	Here is the flower to classify:	Milan criteria:	1907
1857	Input: sepal length: {sepal_length} cm,	- One single tumor with its diameter \leq 5 cm,	1908
1858	sepal width: {sepal_width} cm, petal length:	- Alternatively, up to 3 tumors, each with its diame-	1909
1859	{petal_length} cm, petal width: {petal_width} cm	ter \leq 3 cm.	1910
1860	Output:	Important notes on tumor size interpretation:	1911
1861	Few-shot example in natural language:	- Tumor size (or diameter) may be reported in mil-	1912
1862	Input: sepal length: {sepal_length} cm,	limeters (mm) or centimeters (cm).	1913
1863	sepal width: {sepal_width} cm, petal length:	- It may appear as:	1914
1864	{petal_length} cm, petal width: {petal_width} cm	A single value (e.g., "25 mm"),	1915
1865	Output: {label}	Two dimensions (e.g., "25 x 20 mm"), or	1916
1866		Three dimensions (e.g., "2.5 x 2.0 x 1.8 cm").	1917
1867	Code (Python) prompt.	- In all cases, the diameter is defined as the largest	1918
1868	def classify_iris_species(text: str) -> Lit-	single dimension.	1919
1869	eral["setosa", "versicolor", "virginica"]:	- If a patient has multiple tumors, their sizes will be	1920
1870	"""Classify the species of the iris flower based	listed together as a comma-separated text, e.g., '5	1921
1871	on its features and the following rule.	mm, 30 mm, 3x4x5 mm'.	1922
1872	Rule: If Petal Width $<$ 0.8 cm, then the species	Your task: Read the provided text describing the	1923
1873	is "setosa".	patient's tumor size(s), and classify whether the	1924
1874	If Petal Width \geq 0.8 cm and $<$ 1.8 cm, then the	patient meets the Milan criteria.	1925
1875	species is "versicolor".	If the patient is within the criteria, label 'True'. If	1926
		not, label 'False'.	

1927	Please answer directly with 'True' or 'False' only.	text = "The tumor size(s) of the HCC patient:	1979
1928	{few_shot_examples}	{text}"	1980
1929	Here is the text to classify:	assert (HCC_staging(text) == "{label}")	1981
1930	Input: The tumor sizes of the HCC patient: {text}		1982
1931	Output:	Code (JavaScript) prompt.	1983
1932	Few-shot example in natural language:	const assert = require('assert');	1984
1933	Input: The tumor size(s) of the HCC patient: {text}		1985
1934	Output: {label}	function HCCStaging(text) {	1986
1935		/*	1987
1936	Code (Python) prompt.	Classify the patient according to Milan criteria	1988
1937	def HCC_staging(text: str) -> Literal["True",	based on tumor size(s) in the text.	1989
1938	"False"]:		1990
1939	"""Classify the patient according to Milan	Milan criteria:	1991
1940	criteria based on tumor size(s) in the text.	- One single tumor with its diameter <= 5 cm,	1992
1941	Milan criteria:	- Alternatively, up to 3 tumors, each with its	1993
1942	- One single tumor with its diameter <= 5 cm,	diameter <= 3 cm.	1994
1943	- Alternatively, up to 3 tumors, each with its		1995
1944	diameter <= 3 cm.	Important notes on tumor size interpretation:	1996
1945	Important notes on tumor size interpretation:	- Tumor size (or diameter) may be reported in	1997
1946	- Tumor size (or diameter) may be reported in	millimeters (mm) or centimeters (cm).	1998
1947	millimeters (mm) or centimeters (cm).	- It may appear as:	1999
1948	- It may appear as:	* A single value (e.g., "25 mm"),	2000
1949	* A single value (e.g., "25 mm"),	* Two dimensions (e.g., "25 x 20 mm"), or	2001
1950	* Two dimensions (e.g., "25 x 20 mm"), or	* Three dimensions (e.g., "2.5 x 2.0 x 1.8	2002
1951	* Three dimensions (e.g., "2.5 x 2.0 x 1.8	cm").	2003
1952	cm").	- In all cases, the diameter is defined as the	2004
1953	- In all cases, the diameter is defined as the	largest single dimension.	2005
1954	largest single dimension.	- If a patient has multiple tumors, their sizes	2006
1955	- If a patient has multiple tumors, their sizes	will be listed together as a comma-separated text,	2007
1956	will be listed together as a comma-separated text,	e.g., '5 mm, 30 mm, 3x4x5 mm'.	2008
1957	e.g., '5 mm, 30 mm, 3x4x5 mm'.		2009
1958	This function reads the provided 'text' de-	This function reads the provided 'text'	2010
1959	scribing the patient's tumor size(s), and classify	describing the patient's tumor size(s),	2011
1960	whether the patient meets the Milan criteria.	and classify whether the patient meets the Milan	2012
1961	If the patient is within the criteria, return "True".	criteria.	2013
1962	If not, return "False".	If the patient is within the criteria, return "True".	2014
1963	Args:	If not, return "False".	2015
1964	- text (str): The input text of the HCC patient	Args:	2016
1965	reporting tumor size(s).	- text (string): The input text of the HCC	2017
1966	Returns:	patient reporting tumor size(s).	2018
1967	- str: "True" if the patient meets the Milan	Returns:	2019
1968	criteria, "False" otherwise.	- string: "True" if the patient meets the Milan	2020
1969	"""	criteria, "False" otherwise.	2021
1970	pass	*/	2022
1971		// TODO: Implement Milan staging logic	2023
1972	{few_shot_examples}	}	2024
1973	# Test case for inference		2025
1974	text = "The tumor size(s) of the HCC patient:	{few_shot_examples}	2026
1975	{text}"	# Test case for inference	2027
1976	assert (HCC_staging(text) == "	const text = "The tumor size(s) of the HCC patient:	2028
1977	Few-shot example in Python:	{text}";	2029
1978	# Example test case	assert.strictEqual(HCCStaging(text), "	2030

2031	Few-shot example in JavaScript:	# Test case for inference	2083
2032	# Example test case	std::string text = "The tumor size(s) of the HCC	2084
2033	const text = "The tumor size(s) of the HCC patient:	patient: { <i>text</i> ";	2085
2034	{ <i>text</i> ";	assert(HCC_staging(text) == "	2086
2035	assert.strictEqual(HCCStaging(text), "{ <i>label</i> ");	Few-shot example in C++:	2087
2036		# Example test case	2088
2037	Code (C++) prompt.	std::string text = "The tumor size(s) of the HCC	2089
2038	#include <cassert>	patient: { <i>text</i> ";	2090
2039	#include <string>	assert(HCC_staging(text) == "{ <i>label</i> ");	2091
2040			2092
2041	std::string HCC_staging(const std::string& text)	Code (Ruby) prompt.	2093
2042	{	def HCC_staging(text)	2094
2043	/*	# Classify the patient according to Milan criteria	2095
2044	Classify the patient according to Milan criteria	based on tumor size(s) in the text.	2096
2045	based on tumor size(s) in the text.		2097
2046		# Milan criteria:	2098
2047	Milan criteria:	- One single tumor with its diameter <= 5 cm,	2099
2048	- One single tumor with its diameter <= 5 cm,	- Alternatively, up to 3 tumors, each with its	2100
2049	- Alternatively, up to 3 tumors, each with its	diameter <= 3 cm.	2101
2050	diameter <= 3 cm.		2102
2051		# Important notes on tumor size interpreta-	2103
2052	Important notes on tumor size interpretation:	tion:	2104
2053	- Tumor size (or diameter) may be reported in	- Tumor size (or diameter) may be reported in	2105
2054	millimeters (mm) or centimeters (cm).	millimeters (mm) or centimeters (cm).	2106
2055	- It may appear as:	- It may appear as:	2107
2056	* A single value (e.g., "25 mm"),	* A single value (e.g., "25 mm"),	2108
2057	* Two dimensions (e.g., "25 x 20 mm"), or	* Two dimensions (e.g., "25 x 20 mm"), or	2109
2058	* Three dimensions (e.g., "2.5 x 2.0 x 1.8	* Three dimensions (e.g., "2.5 x 2.0 x 1.8	2110
2059	cm").	cm").	2111
2060	- In all cases, the diameter is defined as the	- In all cases, the diameter is defined as the	2112
2061	largest single dimension.	largest single dimension.	2113
2062	- If a patient has multiple tumors, their sizes	- If a patient has multiple tumors, their sizes	2114
2063	will be listed together as a comma-separated text,	will be listed together as a comma-separated text,	2115
2064	e.g., '5 mm, 30 mm, 3x4x5 mm'.	e.g., '5 mm, 30 mm, 3x4x5 mm'.	2116
2065			2117
2066	This function reads the provided 'text'	# This function reads the provided 'text'	2118
2067	describing the patient's tumor size(s),	describing the patient's tumor size(s),	2119
2068	and classify whether the patient meets the Milan	# and classify whether the patient meets the	2120
2069	criteria.	Milan criteria.	2121
2070	If the patient is within the criteria, return "True".	# If the patient is within the criteria, return	2122
2071	If not, return "False".	"True". If not, return "False".	2123
2072	Args:	# Args:	2124
2073	- text (std::string): The input text of the HCC	- text (string): The input text of the HCC	2125
2074	patient reporting tumor size(s).	patient reporting tumor size(s).	2126
2075	Returns:	# Returns:	2127
2076	- std::string: "True" if the patient meets the	- string: "True" if the patient meets the Milan	2128
2077	Milan criteria, "False" otherwise.	criteria, "False" otherwise.	2129
2078	*/	# TODO: Implement Milan staging logic	2130
2079	// TODO: Implement Milan staging logic	end	2131
2080	}	{ <i>few_shot_examples</i> }	2132
2081		# Test case for inference	2133
2082	{ <i>few_shot_examples</i> }		2134

```
2135 text = "The tumor size(s) of the HCC patient:
2136 {text}"
2137 raise "Assertion failed" unless (HCC_staging(text)
2138 == "
2139 Few-shot example in Ruby:
2140 # Example test case
2141 text = "The tumor size(s) of the HCC patient:
2142 {text}"
2143 raise "Assertion failed" unless (HCC_staging(text)
2144 == "{label}")
2145
```