

# SpecDec

Avi Arya & Shreya Singh & Theo Urban  
10-423/623 Generative AI Course Project

November 23, 2025

## 1 Introduction

As large language models have continued to grow, inference latency has become a bigger issue. Improvements in inference latency can significantly affect workflows, and having lower latency may result in faster systems and feedback loops. A recent proposal to address the issue of inference latency is Speculative Decoding, a method of using smaller models to generate a response, then using a larger language model to attempt to either accept or disregard the smaller language model's result. If the larger language model does not accept the smaller language model's result, then the larger language model generates its own result fully. However, the hope is that, oftentimes, the smaller language model and the larger language model will result, resulting in improved latency.

In this project, we hope to apply finetuning to Speculative Decoding to create a faster model for code generation. Specifically, we will use The Stack, a dataset of code files, to finetune minGPT, then use this as the small model in our speculative decoding model, eventually comparing this model to speculative decoding without this additional finetuning.

## 2 Dataset / Task

For this project, we will use The Stack, a large dataset of source code developed by the BigCode project. The dataset is one of the most extensive collections of code, containing 6.4 terabytes spanning 358 programming languages.

Our use of this dataset will be to finetune our draft model, minGPT 2. We intend to create a training split with code from higher resource languages (Python, C++, Javascript, etc) and perhaps some lower resource languages as well. The lower resource languages can be used to evaluate the finetuning process for its impact on cross lingual performance disparities.

The test set will comprise data from The Stack not used during finetuning. The set will have several code completion prompts where we will use a function signature, comment block, or partial line of code as input.

Using this will allow us to compare the base model and finetuned model fairly on their code generation ability.

The primary task at hand is domain-specific accelerated code completion. In other words, given an input context  $s$  (a string of source code), we want to generate a plausible and syntactically correct continuation of the code using a speculative decoding framework to potentially accelerate the rate of output.

We define our two primary models as follows:

1. **Verifier Model** ( $P_\theta$ ) represents the probability distribution  $p(x | s)$  over the next token  $x$  given some context  $s$ . For this project, we will use the standard GPT 2 model as  $P_\theta$ .
2. **Draft Model** ( $Q_\phi$ ) represents the probability distribution  $q(x | s)$ . For this project, we will use minGPT 2 which is smaller and faster than the verifier model.

Our project will compare two distinct draft models:

1. **Baseline Draft Model** ( $Q_\phi$ ) which is the unmodified minGPT 2 model.
2. **Finetuned Draft Model** ( $Q'_\phi$ ) which is the minGPT 2 model after it has been finetuned on the training split from The Stack.

The task involves providing the system with a context string  $s$  and sampling a fixed length sequence of tokens from the draft model. These tokens are then validated by the verifier in a single pass.

To quantitatively evaluate our model's success, we will look at its acceptance rate, speedup over the baseline model, and crosslingual disparity. These metrics are detailed more in a later section.

## 3 Related Work

Reducing inference times in LLMs is a significant area of past and current research. Modern LLMs have high associated computation and memory bandwidth costs due to standard decoding, which in turn provide a bottleneck for many real world applications. Our project

78 aims to build on the core technique of speculative de- 130  
79 coding, first presented in ?. The paper demonstrated 131  
80 that using a smaller draft model to generate tokens and 132  
81 then passing them to a larger verifier model could pro-  
82 vide speedups in generations without sacrificing gen-  
83 eration quality. Then ? introduced a similar method  
84 called speculative sampling that further solidified the  
85 technique’s potential at reducing overhead. Our base-  
86 line model will be a direct implementation of the algo-  
87 rithms presented in these two papers.

88 The reason for our finetuned model, however, is that 133  
89 the benefits of speculative decoding are not uniformly 134  
90 distributed. This issue is raised and explored by ?, who 135  
91 point out performance disparities in speculative decod- 136  
92 ing. They show that it accelerates generation primar- 137  
93 ily for high resource languages and predictable text se- 138  
94 quences. To this extent, we aim to evaluate our model 139  
95 on crosslingual performance as well. We will adopt the 140  
96 disparity metric  $U(T)$  that they introduced to measure 141  
97 our approach’s success.

98 Our central hypothesis is that specializing the draft 142  
99 model to a particular task can improve its performance, 143  
100 which is well supported in recent literature. An exam- 144  
101 ple is using knowledge distillation, which is done by 145  
102 ? via DistillSpec. Here, the draft model is explicitly 146  
103 trained on the verifier’s output distribution. However, 147  
104 our approach diverges slightly to finetune on a special- 148  
105 ized corpus which can be seen as domain-specific dis-  
106 tillation.

107 Finally, another consideration that we have to make 149  
108 is alignment. As shown by ?, it is imperative that the 150  
109 verifier and draft model are aligned, especially when 151  
110 the larger model is finetuned in some way (i.e. for chat) 152  
111 because a generic draft model will perform quite poorly 153  
112 in such cases. The solution they use is to finetune the 154  
113 draft model using the verifier model as a “teacher”. 155  
114 While our project is not specifically aligning the draft 156  
115 model with the verifier, it does aim to align the draft 157  
116 model to the domain we are testing on.

## 117 4 Approach

118 In this project, our overall goal is to determine the im- 158  
119 pact of finetuning on speculative decoding models. Our 159  
120 baseline model is a speculative decoding model that 160  
121 uses minGPT 2 as the smaller model and GPT 2 as 161  
122 our larger, verifier model. We first must implement 162  
123 speculative decoding using these two models, and run 163  
124 experiments to determine the speed at which this gen- 164  
125 erates its output on a few sample prompts. Next, we 165  
126 will finetune minGPT 2 on The Stack, a large dataset 166  
127 from Hugging Face containing code files in 300 differ- 167  
128 ent languages. By using this dataset, we hope to better 168  
129 inform our smaller model on this particular task. We 169

will then reimplement speculative decoding using this 170  
finetuned model alongside GPT2, and will evaluate the 171  
inference latency and results on our given prompts. 172

## 5 Experiments

### 5.1 Experimental Setup

We will conduct all of our experiences on an A100 ob- 173  
tained through Google Colab Pro. We use ???B Pythia 174  
by ElutherAI as our Verifier model and 70M Pythia as 175  
our Draft model, following from the initial experiments 176  
below. All experiments are conducted with a fixed seed 177  
to ensure consistency across runs.

For our main set of experiemnts, we fix  $\gamma = ???$  to- 178  
kens to standardize across runs with a reasonable bal- 179  
ance between draft and verification speed.

### 5.2 Baseline Performance

Our first experiment establishes the baseline perfor- 180  
mance of a few different models and on different tasks 181  
- natural language completions and a few coding tasks 182  
in different languages outlined in the appendix.

| Model                  | Completion  | Acce. Rate | Speedup |
|------------------------|-------------|------------|---------|
| GPT2-Large, distilgpt2 | NL          | ???        | ???     |
|                        | Code (Py)   | ???        |         |
|                        | Code (C)    | ???        |         |
|                        | Code (go)   | ???        |         |
|                        | Code (rust) | ???        |         |
| Qwen2.5 (7B, 0.5B)     | NL          | ???        | ???     |
|                        | Code (Py)   | ???        |         |
|                        | Code (C)    | ???        |         |
|                        | Code (go)   | ???        |         |
|                        | Code (rust) | ???        |         |
| Pythia (12B, 70M)      | NL          | ???        | ???     |
|                        | Code (Py)   | ???        |         |
|                        | Code (C)    | ???        |         |
|                        | Code (go)   | ???        |         |
|                        | Code (rust) | ???        |         |

Table 1: Caption

### 5.3 Expected Outcomes

We want to answer a few fundamental research ques- 183  
tions:

1. How does domain-specific finetuning of a draft 184  
model on code affect the Acceptance rate and 185  
speedup of a Speculative Decoding setup in and 186  
out of domain.

- 156 2. How does finetuning a draft model on a code  
 157 impact performance on low-resource sets (in our  
 158 case, languages), within the domain.  
 159 3. How does dynamically varying draft length ( $\gamma$ )  
 160 impact speedup under a fixed compute environ-  
 161 ment and under different domains

## 162 6 Proposed Experiments

163 In this section, we outline the experiments designed to  
 164 evaluate how domain-specific finetuning of the draft  
 165 model influences speculative decoding effectiveness.  
 166 Our proposed experiments address three research ques-  
 167 tions: (1) whether finetuning improves acceptance rate  
 168 and speedup in- and out-of-domain, (2) how per-  
 169 formance changes across high- and low-resource lan-  
 170 guages, and (3) how the draft length  $\gamma$  affects perfor-  
 171 mance under a fixed compute budget.

### 172 6.1 RQ1: Does Domain-Specific Finetun- 173 ing Improve Speculative Decoding Per- 174 formance?

175 **Goal.** Quantify the effect of finetuning on acceptance  
 rate, latency, and speedup across code and natural lan-  
 guage tasks.

**Experimental Setup.** We compare two systems:

1. **Baseline SpecDec:** GPT2 (verifier) + unfinetuned minGPT2 (draft)
2. **Finetuned SpecDec:** GPT2 (verifier) + minGPT2 finetuned on The Stack

Both models will be evaluated on:

- **In-domain:** Code completions across Python, C++, Go, Rust, and others
- **Out-of-domain:** Natural language tasks (news, stories, instructions)

Each prompt set will generate 128-token continuations with fixed draft length  $\gamma$ .

| Model Pair | Domain    | Accept Rate | Latency | Speedup | $\gamma \in \{1, 2, 4, 8, 16, 32\}$ |
|------------|-----------|-------------|---------|---------|-------------------------------------|
| Baseline   | Code (Py) | ???         | ???     | ???     |                                     |
| Finetuned  | Code (Py) | ???         | ???     | ???     | For each $\gamma$ , we record:      |
| Baseline   | NL        | ???         | ???     | ???     |                                     |
| Finetuned  | NL        | ???         | ???     | ???     | • Acceptance rate                   |

Table 2: Effect of draft model finetuning on acceptance rate and decoding speed.

**Expected Outcome.** Finetuning should increase ac-  
 ceptance rate and speedup on code tasks, while main-  
 taining comparable performance on natural language  
 prompts.

### 6.2 RQ2: Does Finetuning Improve Perfor- mance for Low-Resource Languages? 194 195

**Goal.** Determine whether finetuning disproportio-  
 nately helps high-resource languages or narrows cross-  
 lingual disparities. 196  
 197  
 198

**Experimental Setup.** We construct a language-  
 stratified test set from The Stack: 199  
 200

- **High-resource:** Python, C++, Java 201
- **Medium:** Lua, Haskell, Ruby 202
- **Low-resource:** Zig, Elixir, OCaml 203

For each language bucket, we measure: 204

- Acceptance rate 205
- End-to-end speedup 206
- Disparity metric  $U(T)$  from ? 207

| Language Group | Baseline $U(T)$ | Finetuned $U(T)$ | $\Delta$ Dispary |
|----------------|-----------------|------------------|------------------|
| High Resource  | ???             | ???              | ???              |
| Medium         | ???             | ???              | ???              |
| Low Resource   | ???             | ???              | ???              |

Table 3: Change in cross-lingual performance disparity  
 after finetuning.

**Expected Outcome.** Finetuning should reduce per-  
 formance disparity across languages by improving  
 draft alignment with code structure. 208  
 209  
 210

### 6.3 RQ3: How Does Draft Length $\gamma$ Affect Speedup Under Fixed Compute? 211 212

**Goal.** Characterize how varying  $\gamma$  impacts accept-  
 ance rate, rejection bursts, and overall speedup. 213  
 214

**Experimental Setup.** Using the finetuned draft  
 model, we vary: 215  
 216

- Rejection burst statistics 219
- Real-time decoding speed 220

**Expected Outcome.** We hypothesize a unimodal  
 trend where moderate values of  $\gamma$  maximize speedup  
 without causing excessive rejection cascades. 221  
 222  
 223

| $\gamma$ | Avg. Burst Len. | $P(\text{Burst} > 10)$ | Speedup | Accept Rate |
|----------|-----------------|------------------------|---------|-------------|
| 1        | ???             | ???                    | ???     | ???         |
| 4        | ???             | ???                    | ???     | ???         |
| 16       | ???             | ???                    | ???     | ???         |

Table 4: Effect of draft length on rejection patterns and speedup.

## 6.4 Summary

These experiments collectively evaluate:

- Finetuning effects on speculative decoding efficiency (RQ1) 224
- Cross-lingual disparity and fairness implications (RQ2) 225
- Sensitivity of speculative decoding to draft length (RQ3) 226

Together, they form a comprehensive framework for understanding speculative decoding in domain-specialized settings. 227

## 7 Compute

## 8 Plan

We can break this project in to three sequential steps:

1. Set up speculative decoding using minGPT and GPT2 as paired models 232
2. Create metrics 233
3. Finetune minGPT and GPT2 on The Stack 234
4. Evaluate metrics pre and post finetune 235

We intend to complete through the metrics creation and initial testing by the executive summary deadline. The work is highly sequential so it will involve a lot of pair programming and back-and-forth. However, we will assign directly responsible members for items 1, 2, and 3, to Avi, Shreya, and Theo in that order. 236