JSONSchemaBench: A Rigorous Benchmark of Structured Outputs for LLMs

Anonymous ACL submission

Abstract

Reliably generating structured outputs has become a critical capability for modern language model (LM) applications. Constrained decoding has emerged as the dominant technology across sectors for enforcing structured outputs during generation. Despite its growing adoption, little has been done with the systematic evaluation of the behaviors and performance of constrained decoding. Constrained decoding frameworks have standardized around JSON Schema as a structured data format, with most uses guaranteeing constraint compliance given a schema. However, there is poor understanding of the effectiveness of the methods in practice. We present an evaluation framework to assess constrained decoding approaches across three critical dimensions: efficiency in generating constraint-compliant outputs, coverage of diverse constraint types, and quality of the generated outputs. To facilitate this evaluation, we introduce JSONSchemaBench, a benchmark for constrained decoding comprising 10K real-world JSON schemas that encompass a wide range of constraints with varying complexity. We find that JSONSchemaBench presents a significant challenge for both LLMs and constrained decoding frameworks, highlighting ample room for improvement and exposing gaps in the existing solutions.

1 Introduction

007

017

031

042

The rapid advancements in LMs in recent years have significantly broadened their applications, extending beyond natural language tasks to more complex challenges such as web navigation (Yao et al., 2023b), data extraction (Polak and Morgan, 2024), and tool use (Schick et al., 2023). Unlike traditional natural language processing (NLP) tasks where the output is aimed at review by humans, output in these applications is often consumed by machines such as controller and service APIs. The machineoriented nature of these applications requires LMs



Figure 1: Comparison across various constraineddecoding frameworks by efficiency (speed of output generation), coverage (support for JSON Schema features), and quality (effects on underlying task accuracy). Guidance outperforms other frameworks on these dimensions.

to generate structured outputs that strictly adhere to predefined formats and constraints. However, the LM generation process is probabilistic and does not provide guarantees on the output's structure, making it challenging to deploy LMs in applications requiring structured inputs and high reliability.

The methodology of constrained decoding, a technique that integrates constraints into the decoding process of LMs, has been developed to address the need to adapt LM generations to the challenge of providing structured output. Constrained decoding intervenes in the decoding process of LMs by masking out invalid tokens based on given constraints and prefix tokens. This intervention guides the LM to sample only from valid tokens, ensuring that the final output perfectly conforms to a predefined structure.

The strong demand for structured generation (Liu et al., 2024) has led to the development of

061

043

various constrained-decoding frameworks¹, such as Guidance (Guidance AI, 2023), Outlines (Willard and Louf, 2023), XGrammar (Dong et al., 2024) and the grammar module of Llamacpp (Gerganov and al., 2023) These frameworks provide broad support for different types of constraints, minimal overhead, and compatibility with various LM ecosystems, facilitating the adoption of constrained decoding in real-world applications.

062

063

064

067

072

086

092

097

100

101

103

104

105

106

107

108

109

110

JSON Schema offers a high level, domainspecific way to define constraints for JSON data, a widely adopted data interchange format. As a result, JSON Schema has emerged as a key specification language for constrained decoding. Commercial LM providers, such as OpenAI, have embraced constrained decoding by incorporating support for JSON Schema directly into their APIs. These integrations highlight the emergence of JSON Schema as an industry-wide standard for specifying constraints on structured outputs, ensuring compatibility across diverse applications. Despite the growing adoption of constrained decoding for structured generation, several issues and questions persist:

Q1: Efficiency: Does constrained decoding slow down or speed up the generation process? Which framework is the most efficient?

Q2: Coverage: The JSON Schema specification has an evolving and expansive feature set. How well do existing constrained decoding frameworks support these features?

Q3: Quality: While constrained decoding guarantees that LM outputs conform to a desired structure, does it negatively affect the semantic quality of outputs?

To answer these questions, we need to study constrained-decoding methods with a large-scale, diverse, and real-world collection of user-defined structures. To evaluate the performance of constrained decoding frameworks, we introduce JSONSchemaBench, a collection of 10K realworld JSON schemas from various sources, Organized into 10 datasets of varying complexity and diversity, the benchmark spans domains such as function signatures, service APIs, and system configurations. We evaluate six state-of-the-art constrained decoding frameworks, including Guidance, Outlines, Llamacpp, XGrammar, OpenAI, and Gemini, on JSONSchemaBench. We pair this real-world schema dataset with the official JSON Schema Test

Suite (JSON Schema Org, 2024) in order to extract 111 detailed insights into coverage of JSON Schema 112 functionality across these frameworks, and to fur-113 ther evaluate them with considerations of end-to-114 end task accuracy in the context of multiple real-115 world tasks. Altogether, our evaluation takes three 116 aspects into consideration: efficiency, coverage, 117 and quality. We define specific metrics to mea-118 sure these three functional aspects and evaluate 119 constrained decoding frameworks against them. 120 Through extensive experiments, we converge on 121 the following findings as illustrated in Figure 1. (1) 122 Constrained decoding can speed up the generation 123 process by 50% compared to unconstrained decod-124 ing. (2) Frameworks demonstrate significant differ-125 ences in their actual support for real-world JSON 126 schemas, with the best framework supporting twice 127 as many schemas as the worst. (3) Constrained 128 decoding consistently improves the performance 129 of downstream tasks up to 4%, even for tasks with 130 minimal structure like GSM8k. 131

Contributions Our contributions are three-fold:

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

- We assemble JSON schemas from various sources and organize them into a benchmark, JSONSchemaBench, to facilitate the evaluation of constrained decoding frameworks on JSON schema.
- We propose a fine-grained evaluation framework to assess the versatility of constrained decoding frameworks in handling diverse JSON schema features, including declared coverage, empirical coverage, and compliance rate.
- We evaluate six state-of-the-art constrained decoding frameworks on JSONSchemaBench, uncovering their strengths and limitations in generating schema-compliant JSON outputs and analyzing their impact on downstream tasks.

2 Background and Related Work

JSON Schema is a meta-language that describes the structure of JSON data. It is capable of expressing a wide variety of constraints, such as the types of JSON object properties, the length of JSON arrays or the pattern that a JSON string must match. The syntax and capabilities of JSON Schema are defined in the JSON Schema specification (Wright et al., 2022), which defines a large number of *keywords*, each of which may be used or combined with other keywords within a schema to enforce

¹We use the terms *constrained decoding framework* and *grammar engine* interchangeably.

188

160

161

162

163

164

165

166

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

185

186

187

193

194

195 196 constraints like the ones mentioned. JSON Schema is widely used in the software ecosystem, and previous work has been done to collect extensive examples of JSON Schemas with a focus both on real-world use as well as on overall correctness.

Baazizi et al. (2021) collected over 6,000 JSON schemas from publicly available GitHub repositories. Attouche et al. (2022) used it alongside additional collected JSON schemas in order to evaluate a witness generation algorithm for JSON Schema. Separately, the official JSON Schema Test Suite (JSON Schema Org, 2024) is a collection of manually created test cases, maintained by the JSON Schema core team, which exercises a large portion of the functionality defined in the JSON Schema specification. It was originally written to assist implementers of JSON Schema validation tools with testing their compliance against the specification, and therefore contains a wide variety of examples for each of JSON Schema's keywords, including in edge case scenarios. Notably, Bowtie (Bowtie, 2025) leverages the test suite as a foundation for comparing and understanding different implementations of the JSON Schema specification across programming languages. Taken together, these two datasets form a large number of examples both of JSON Schema's diverse feature set as well as its use in the wild.

Algo	Algorithm 1 Constrained Decoding			
Requ	ire: Constraint C, LLM f	, Prompt x		
Ensu	re: Output o adhering to C	γ		
1: c	$\rightarrow \leftarrow [] \rightarrow \bullet$			
2: I	oop			
3:	C.update(o)	\triangleright advance state of C		
4:	$m \leftarrow C.mask()$	⊳ compute mask		
5:	$v \leftarrow f(x+o)$	⊳ compute logits		
6:	$v' \leftarrow m \odot v'$			
7:	$t \leftarrow decode(\alpha')$	⊳ sample		
8:	if $t = EOS$ then	-		
9:	break			
10:	end if			
11:	o.append(t)			
12: 6	end loop			
13: i	eturn o	⊳ output		

Constrained decoding (Deutsch et al., 2019; Shin et al., 2021; Scholak et al., 2021; Poesia et al., 2022; Wang et al.; Geng et al., 2023) refers to methods that guide the generation process of language models (LMs) by masking out tokens that do not adhere to predefined constraints at each step. Recently, highly optimized grammar-constrained decoding frameworks (Guidance AI, 2023; Beurer-Kellner et al., 2023; Willard and Louf, 2023; Kuchnik et al., 2023; Zheng et al., 2024; Dong et al., 2024) have been developed to improve the efficiency and usability of constrained decoding. 197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

228

229

230

232

233

234

235

236

237

238

239

240

241

242

243

The evaluation of constrained decoding remains an under-explored topic, with no consensus on what defines the effectiveness of constrained decoding. While some research has pursued comparisons of constrained decoding with unconstrained LMs (Roy et al., 2024; Tang et al., 2024; Yao et al., 2023a), the studies to date fail to provide comparisons across different constrained decoding frameworks. The benchmarks employed have either narrowly focused on specific tasks or rely on formalgrammar–based artificial setups, that have unclear relevance to real-world use cases.

3 The JSONSchemaBench

Our goal is to design a benchmark that is (1) diverse enough to cover the most common constraint types encountered in real-world applications, (2) large enough to provide a reliable evaluation, and (3) equipped with fair and multidimensional metrics to ensure comprehensive assessments.

3.1 Data Collection

We start with the 6K JSON schemas collected by (Baazizi et al., 2021) from publicly available GitHub repositories, and with the set of schemas from the JSON Schema Test Suite (JSON Schema Org, 2024). We further collect JSON schemas from other sources, such as the JSON Schema Store (Schema Store Org, 2014), the GlaiveAI function calling dataset V2 (GlaiveAI, 2024), and from Kubernetes configuration files (Kubernetes, 2022). We filter out invalid schemas and standardize the schemas to ensure that they conform to the version of JSON Schema declared² in each schema The GitHub JSON schemas collection from (Baazizi et al., 2021) contains schemas of varying complexity and diversity, ranging from simple type constraints to complex constraints with nested objects and arrays. For more fine-grained evaluation, we split the data into five collections based on the schema size: trivial, small, medium, large, ultra. The suites finalized after all collection and processing are listed in Table 1. We excluded GitHub-Trivial and GitHub-Ultra from the experiments as they were too easy or too hard. However, we retained these datasets in the benchmark, with

²The \$schema keyword, defined in the JSON Schema specification, allows any schema to self-identify which version of JSON Schema it is written for.

Dataset	Category	Count
GlaiveAI-2K	Function Call	1707
Github-Trivial	Misc	444
Github-Easy	Misc	1943
Snowplow	Operational API	403
Github-Medium	Misc	1976
Kubernetes	Kubernetes API	1064
Washington Post	Resource Access API	125
Github-Hard	Misc	1240
JSONSchemaStore	Misc	492
Github-Ultra	Misc	164
Total		9558

Table 1: Schema collection metadata.

GitHub-Ultra serving as an aspirational target for future advancements. For more information on post-processing and dataset splitting, we refer the reader to Appendix A.

4 Efficiency

244

245

247

249

251

255 256

257

258

261

262

263

265

267

268

270

271

272

Naïve implementations of constrained decoding add overhead to the standard LM inference process, including a per-step mask computation and an optional one-time grammar compilation. However, several optimizations can significantly reduce this overhead. For instance, mask computation can run in parallel with the LM's forward pass, and grammar compilation can be performed concurrently with pre-filling computations (Guidance AI, 2023; Dong et al., 2024). Other optimizations such as grammar caching and constraint-based speculative decoding (GuidanceAI, 2024a; Beurer-Kellner et al., 2023; Kurt, 2024a) can further reduce overhead.

Metrics We break down the efficiency evaluation into the following components:

- Grammar Compilation Time (GCT): The time spent on grammar compilation.
- **Time to First Token (TTFT):** Time from the start of generation to the production of the first token.
- **Time per Output Token (TPOT):** Average time to generate each token after the first.

4.1 Setup

The efficiency experiment depends on both the
size of the model and the tokenizer's vocabulary
size. We used Llama-3.1-8B-Instruct with the
Llamacpp inference engine as backend for Outlines, Guidance, and Llamacpp. As XGrammar
doesn't support Llamacpp as backend, we add

an additional experiment with the **Hugging Face Transformers** inference engine for XGrammar. All experiments are conducted on a single **NVIDIA A100-SXM4-80GB GPU** with **AMD EPYC 7543** (**12 cores**) CPU. The batch size is set to 1 for all experiments. Additional details about setup are provided in the Appendix E. We also provide a snippet of how we call each engine in the Appendix G.

Addressing coverage bias. The efficiency metrics are meaningful only for instances that a grammar engine can process. Different engines exhibit varying levels of schema coverage, with some engines handling a wider range of schemas than others. Engines with lower coverage often process simpler, shorter schemas, which naturally compile and generate faster. As a result, averaging efficiency metrics across covered instances can introduce bias favoring engines with lower coverage. For a more detailed discussion on coverage, see Section 5. To ensure fairness, we calculate efficiency metrics on the intersection of covered instances across all engines.

Dataset	Framework	GCT	TTFT	ТРОТ
		(s)	(s)	(ms)
GlaiveAI	LM only	NA	0.10	15.40
	Guidance	<u>0.00</u>	0.24	<u>6.37</u>
	Llamacpp	0.05	0.20	29.98
	Outlines	3.48	3.65	30.33
GitHub	LM only	NA	0.10	15.83
Easy	Guidance	0.00	$\overline{0.34}$	7.44
•	Llamacpp	$\overline{0.05}$	0.18	$2\overline{7.22}$
	Outlines	3.71	3.97	39.78
Snowplow	LM only	NA	0.11	16.23
	Guidance	0.00	0.28	6.55
	Llamacpp	0.05	0.20	$2\overline{8.90}$
	Outlines	3.91	4.14	42.66
GitHub	LM only	NA	0.20	16.68
Medium	Guidance	<u>0.01</u>	0.54	7.57
	Llamacpp	0.06	0.30	29.08
	Outlines	8.05	8.38	46.57
Kubernetes	LM only	NA	0.16	15.32
	Guidance	0.01	0.45	<u>9.47</u>
	Llamacpp	0.05	0.28	28.04
	Outlines	5.29	5.55	46.10

Table 2: Efficiency metrics for different engines with LlamaCpp as the inference engine. GCT: Grammar Compilation Time, TTFT: Time to First Token, TPOT: Time Per Output Token. Bold values indicate the smallest in each column for GCT, TTFT, and TPOT. All values are median of the samples. Results for the GitHub Hard and Washington Post datasets are provided in Appendix E.

Grammar compilation time. There are notable differences in grammar compilation times between

279

281

283

284

285

287

290

291

292

294

295

297

298

299

355

356

357

360

361

362

363

364

365

366

367

369

371

372

373

374

375

376

377

378

334

335

336

337

321

303

304

305

307

311

312

324 325 326 327

323

328

3

331

332 333 the engines. Both Guidance and Llamacpp dynamically compute their constraints during token generation, leading to minimal grammar compilation time. In the middle, XGrammar does include a nontrivial compilation step, but they are able to largely mitigate its impact by running it concurrently with prompt pre-filling. Finally Outlines, which converts JSON schemas into regular-expression based constraints, has significantly higher compilation time.

313Time per output token. While Outlines and Lla-314macpp demonstrate substantially lower throughput315than the LM-only approach, Guidance achieves316even higher efficiency, which it accomplishes by317fast-forwarding ³ certain generation steps with318its guidance acceleration (GuidanceAI, 2024a).319Comparing Guidance and XGrammar with the HF320Transformers backend shows that Guidance has a
significantly better TPOT.

Dataset	Framework	GCT (s)	TTFT (s)	TPOT (ms)
GlaiveAI	Guidance	<u>0.01</u>	0.36	<u>36.92</u>
	XGrammar	0.12	<u>0.30</u>	66.78
GitHub	Guidance	<u>0.01</u>	0.37	<u>42.03</u>
Easy	XGrammar	0.11	0.33	65.57
GitHub	Guidance	<u>0.01</u>	0.55	<u>44.21</u>
Medium	XGrammar	0.20	<u>0.48</u>	65.51
GitHub	Guidance	<u>0.01</u>	0.73	<u>35.88</u>
Hard	XGrammar	0.30	0.65	65.20

Table 3: As XGrammar doesn't support **llama.cpp**, we add an additional experiment with the **Hugging Face Transformers** inference engine for XGrammar and Guidance. All values are **median** of the result samples.

5 Coverage

Each constrained decoding framework has limitations when it comes to translating JSON schemas into a set of constraints that can reliably guarantee the validity of LM outputs. To systematically evaluate the effectiveness of these frameworks, we define three notions of coverage:

Definition 5.1 (Declared Coverage) A schema is considered declared covered if the framework processes the schema without explicitly rejecting it or encountering runtime errors such as exceptions or crashes. **Definition 5.2 (Empirical Coverage)** A schema is considered empirically covered if our experiments show that the constraints generated by the framework result in LM outputs that are schemacompliant.

Definition 5.3 (True Coverage) A schema is considered truly covered if the framework produces constraints that are precisely equivalent to the original JSON Schema definition, i.e., permitting all schema-compliant generations while rejecting all schema-noncompliant generations.

The most ideal coverage metric is the true cover*age*, denoted as C_{True} . However, due to the infinite number of JSON instances that could be validated against a schema, it is difficult to measure in practice without a formal verification method that is capable of exhaustively comparing the schema's semantics against the framework's implementation. $C_{\text{Empirical}}$ is an approximation of C_{True} as it only checks whether the finitely many outputs seen during our experiments conform to a given schema⁴. While C_{Declared} is not an estimate of C_{True} per se, it is an upper-bound of both $\mathcal{C}_{\text{Empirical}}$ and C_{True} and is useful in deriving an additional metric from the coverage evaluation: Compliance Rate $= C_{\text{Empirical}}/C_{\text{Declared}}$. The Compliance Rate estimates the reliability of the constrained decoding framework in guaranteeing compliance given it accepts a given schema.

5.1 Setup

To measure empirical coverage, we conduct all experiments using the Llama-3.2-1B-Instruct model as it is small enough to run efficiently while still producing high-quality outputs. The prompt consists of a simple instruction with two-shot examples (Figure 3), and validation is performed using the jsonschema Python library (Berman, 2025) (using JSON Schema Draft2020-12) with string-format checks enabled. We use greedy decoding with zero-temperature, performing a single generation run, and enforce a 40-second timeout for grammar compilation and an additional 40 seconds for generation. Exceeding these limits is treated as a schema processing failure. Additional details are provided in Appendix B.

 $^{^{3}}$ See Tables12 and 13 for the number of tokens fast-forwarded.

⁴Additionally, we define *theoretical coverage* as the proportion of schemas whose features are fully supported by the grammar engine, with details provided in Appendix C.

5.2 Results

379

400

401

402

403 404

405

406

407 408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

497

428

Empirical Coverage Guidance shows the highest empirical coverage on six out of the eight datasets, with Llamacpp taking the lead on the remaining two: the domain-specific Washington Post and notably hard JSON Schema Store. On the other hand, closed-source grammar engines consistently have the lowest coverage; they came in last on all but one dataset. We note that while empirical coverage is a reasonable indicator of a framework's real-world performance, it is influenced by factors such as the LM being used and the sampling methods employed.

Compliance Rate Among open-source engines, guidance consistently demonstrates the highest compliance rate across all datasets, making it the most reliable option for ensuring schema compliance. Outlines has a comparatively lower compliance rate, primarily due to timeouts during generation. Our analysis reveals that JSON Schema features like minItems, maxItems, enum, and Array, while supported, often take 40 seconds to 10 minutes for Outlines to process. While closed-source implementations have low empirical coverage, they have very high compliance rates, indicating that their providers have taken a more conservative strategy, implementing only a subset of JSON Schema features that they can reliably support.

5.3 JSON Schema Test Suite: Complementary Evaluation

Originally designed to test the correctness and compliance of JSON Schema validation implementations, the official JSON Schema Test Suite (JSON Schema Org, 2024) is a comprehensive collection of test cases spanning the many features of the JSON Schema specification. We believe that the test suite is an ideal tool for assessing the correctness of grammar engines.

The test suite organizes its test cases into 45 categories, each of which corresponds to a feature of JSON Schema, typically a specific keyword such as required or group of tightly related keywords such as if-then-else. A small number of additional categories test broader behaviors, such as infinite-loop-detection. Each test case contains a single schema paired with a collection of JSON instances that are marked as either valid or invalid under that schema. For the purpose of evaluating coverage, we assert that an engine must successfully generate each valid instance and block

Dataset	Framework	Declared	Empirical	Compliant
-				Kate
GlaiveAI	Guidance	0.98	<u>0.96</u>	<u>0.98</u>
	Llamacpp	0.98	0.95	0.97
	Outlines	0.99	0.95	0.96
	XGrammar	1.00	0.93	0.93
	OpenAI	0.89	0.89	1.00
	Gemini	0.86	0.86	1.00
GitHub	Guidance	0.90	0.86	0.96
Easy	Llamacpp	0.85	$\overline{0.75}$	$\overline{0.88}$
v	Outlines	0.86	0.59	0.83
	XGrammar	0.91	0.79	0.87
	OpenAI	0.30	0.29	0.97
	Gemini	0.08	0.07	0.88
Snownlow*	Guidance	0.87	0.82	0 94
Showpiow	Llomoonn	0.07	$\frac{0.02}{0.74}$	$\frac{0.74}{0.81}$
	Outlines	0.92	0.74	0.61
	XGrammar	0.95 NA	0.50 NA	0.01 NA
		0.21	0.01	1.00
	OpenAl	0.21	0.21	1.00
GitHub	Guidance	0.79	<u>0.69</u>	<u>0.87</u>
Medium*	Llamacpp	0.77	0.57	0.74
	Outlines	0.72	0.29	0.40
	XGrammar	0.79	0.52	0.66
	OpenAI	0.13	0.12	0.92
Kubernetes*	Guidance	0.98	0.91	0.92
	Llamacpp	0.98	0.76	0.78
	Outlines	0.98	0.57	0.58
	XGrammar	0.12	0.07	0.58
	OpenAI	0.21	0.21	1.00
Washington	Guidance	0.86	0.86	1.00
Post*	Llamacon	0.97	0.94	0.97
	Outlines	0.97	$\frac{1}{0.22}$	0.23
	XGrammar	0.85	0.64	0.75
	OpenAI	0.13	0.13	1.00
GitHub	Guidance	0.60	0.41	0.69
Hard*	Llamacpp	0.61	$\overline{0.39}$	0.63
	Outlines	0.47	0.03	0.06
	XGrammar	0.69	0.28	0.41
	OpenAI	0.09	0.09	1.00
IsonSchome	Guidance	0.35	0.30	0.88
Storo*	Llamaonr	0.55	0.30	0.00
Store	Outline	0.34	0.00	0.09
	VC	0.38	0.09	0.24
	AGrammar	0.76	0.33	0.43
	OpenAI	0.06	0.06	1.00

Table 4: **Coverage of all the frameworks** on JSON-SchemaBench. Empirical coverage between Open Source engines and OpenAI/Gemini are not directly comparable due to differences in the underlying model (Llama 3.2-1B vs. proprietary models).

 * Gemini results are ommitted for dataset suites with < 1% support.

generation of each invalid instance to "pass" a test case. In addition to compilation failures, we define two failure modes that a grammar engine can exhibit:

Definition 5.4 (Over-constrained) A framework is over-constrained if it rejects JSON instances that are vsalid according to a given JSON Schema. This means the engine is too strict and excludes outputs

Framework	Compile-	Over-	Under-
	Error	constrained	constrained
Outlines	42	16	8
Llamacpp	37	18	7
XGrammar	<u>3</u>	<u>5</u>	38
Guidance	25	7	1

Table 5: Number of categories for which each failure type occurred at least once. Rows do not necessarily sum to the total number of categories, as some categories may have more than one failure type or no failures at all. Bold numbers indicate the framework with the fewest number of failures of a given type.

that should be allowed.

Definition 5.5 (Under-constrained) A framework is under-constrained if it allows JSON instances that are invalid according to a given JSON Schema. This means the engine is overly permissive and allows outputs that should be rejected.

An illustration is given in Figure 5 in Appendix D. *Over-constrained* grammar engines risk limiting the expressive power of LMs, potentially preventing the generation of valid responses and negatively impacting downstream task performance. Conversely, under-constrained engines cannot guarantee that all responses will be valid, often necessitating additional post-processing or retry logic.

5.3.1 Results

Coverage Analysis For each grammar engine and category in the test suite, we calculate *test coverage* as the proportion of passing test cases, reported in Figure 6 in Appendix D Additionally, Table 6 aggregates these metrics, counting categories with minimal coverage (> 0%), partial coverage (> 25%), moderate coverage (> 50%), high coverage (> 75%), and full coverage (100%). We indicate the number of categories for which each framework achieves the highest test coverage (either as the single highest or as the sole leader) as well as the number of categories for which each framework is the sole leader.

• Overall Performance: Guidance outperforms other engines at all coverage levels, achieving full coverage on 13 categories and moderate coverage on 21. In comparison, Llamacpp and XGrammar have full coverage on only one category and moderate coverage on five and three categories, respectively, while Outlines has no full coverage on any category and moderate coverage on two categories.

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

• **Single Highest**: Guidance has the single highest coverage in 19 categories, followed by XGrammar with 10, and Outlines with one, and Llamacpp with none.

Failure Analysis Table 5 provides a breakdown of failure modes for each framework across the test suite, detailing the number of categories with compilation errors, failures to generate positive instances (over-constrained), and failures to block negative instances (under-constrained). Overall, Guidance demonstrates the fewest total failures, in particular minimizing under-constrained errors. Outlines, Llamacpp, and Guidance follow a consistent failure pattern, with most errors occurring during compilation and over-constrained failures being more frequent than under-constrained ones. In contrast, XGrammar minimizes compilation errors but shows the highest number of under-constrained failures, indicating a trade-off favoring permissiveness.

We acknowledge that there is no straightforward correspondence between test suite performance and empirical coverage. One reason for this is that not all features are equally represented in real-world schemas. As a result, strong or weak performance on specific features can have disproportionate impacts depending on their prevalence. Another reason is under-constraining effectively delegates responsibility to the LM, which may produce valid output despite a lack of strict constraints. We emphasize that while under-constraining can be a legitimate strategy, it requires careful implementation and transparency to ensure reliability.

6 Quality

In principle, constrained decoding should not affect the quality of the generated output as it only filters out the invalid tokens. However, things become more complicated due to ambiguity of tokenization (Vivien, 2024; GuidanceAI, 2024b; Geng et al., 2024) and the distributional shifts caused by the intervention (Geng et al., 2023; Tam et al., 2024). As a hypothetical toy example, an LM might answer 89,000 instead of the correct 89000 in a GSM8K question. Constrained decoding can block the invalid token , enforcing structural compliance but potentially may cause the LM to go out of distribution and generate 890000 instead. (Kurt,

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467 468

469

470

471

Coverage	Outlines	Llamacpp	XGrammar	Guidance
Minimal coverage (>0%) Partial coverage (>25%) Moderate coverage (>50%) High coverage (>75%) Full coverage (100%)	20 11 2 0 0	21 11 5 2 1	28 16 3 1 1	$\frac{30}{25}$ $\frac{21}{17}$ $\frac{17}{13}$
Tied for highest (>0%) Single highest	4 1	6 0	14 10	$\frac{25}{19}$

Table 6: Number of categories with a given level of coverage. Each row represents a cumulative coverage threshold, with higher thresholds indicating stricter levels of success. Bold numbers indicate the framework with the highest value in that row.

2024b) argued that the performance decline observed in previous studies (Tam et al., 2024) comes from inadequate prompting, insufficient contextual information, and poorly crafted schemas.

6.1 Setup

522

523

524

525

527

528

530

532

534

539

540

541

542

543

544

545

546

547

548

549

551

Kurt (2024b); Tam et al. (2024) have introduced a series of tasks to investigate potential quality concerns in constrained decoding, which we leverage and extend in this benchmark. Specifically, we adopt the three reasoning tasks from these studies to evaluate the impact of constrained decoding on task accuracy, as detailed in Table 14. The simple output structure of these tasks was designed to isolate the effects of constrained decoding on reasoning, as outlined by (Tam et al., 2024). For our experiments, we use the Llama-3.1-8B-Instruct model to measure task performance. We follow the original setup and prompt specifications from (Kurt, 2024b), with full details provided in Appendix F.

We implement the following constraints for the first three tasks: (1) Last Letter the output needs to be a concatenation of letters from a-z; (2) Shuffle Objects the output needs to be a single letter from A-E enclosed in parentheses; (3) GSM8K the output is an valid integer or float number. The outputs for all three tasks are structured as JSON objects with two fields: "reasoning" and "answer", formatted as {"reasoning": <reasoning about the answer>, "answer": <final answer>}.

6.2 Results

552The results in Table 7 show that the constrained553decoding, regardless of the framework, achieves554higher performance than the unconstrained setting.555Among the frameworks evaluated, Guidance con-556sistently delivers the best performance across all557tasks, with approximately a 3% improvement over558the LM-only approach in every task. We believe

this may be attributed to its token-healing implementation (GuidanceAI, 2024b).

	Last Letters	Shuffle Objects	GSM8K
LM only	50.7%	52.6%	80.1%
XGrammar	51.2%	52.7%	83.7%
Llamacpp	52.0%	52.6%	82.4%
Outlines	53.3%	53.0%	81.6%
Guidance	<u>54.0%</u>	<u>55.9%</u>	<u>83.8</u> %

Table 7: Accuracy on the quality tasks.

7 Conclusion

We have proposed a comprehensive evaluation framework for constrained decoding frameworks with JSON schemas, focusing on efficiency, coverage, and output quality. We introduced **JSON-SchemaBench**, a benchmark comprising 10K realworld JSON schemas, to enable robust assessment under realistic conditions. Our evaluation highlights both the advancements and limitations of current state-of-the-art constrained decoding frameworks. We hope our findings and benchmark guide future research in structured generation, helping the community identify effective tools and extend capabilities with constrained decoding.

8 Limitations

1. Short-Circuited Generation In the experiments, the LM sometimes produces minimal, schema-compliant JSON objects that are trivially correct but lack semantic depth. On most datasets, the average output size is 50 tokens, occasionally as low as 20 tokens (e.g., Kubernetes' and Washington Post'), far below the expected 100 tokens. Such outputs often contain only a few key-value pairs or even empty dictionaries, might inflating the evaluation of constraint compliance. To address this, we propose: 561

562

563

564

566

567

568

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

587

- 593
- 594
- 596

611 612

613

614 615 616

617 618

621 622

630 631

633

637

- **Contextual Prompts:** Enrich prompts with contextual information to guide the model in generating more meaningful outputs.
- Longer Few-Shot Examples: Include longer, more detailed examples as few-shot to encourage generation of long and non-trivial JSON objects.

2. Lack of Semantic Evaluation Our current evaluation emphasizes schema compliance but overlooks the semantic accuracy of the generated JSON objects. While output quality is assessed in Section 6, the tasks primarily involve reasoning tasks with minimal structure. Future work should include semantic evaluation in more complex tasks, such as function calls within LLM agents, to provide a holistic understanding of output quality.

References

- Snowplow Analytics. 2022. Iglu central. https: //github.com/snowplow/iglucentral. Commit hash 726168e. Retrieved 19 September 2022.
- Lyes Attouche, Mohamed-Amine Baazizi, Dario Colazzo, Giorgio Ghelli, Carlo Sartiani, and Stefanie Scherzinger. 2022. Witness Generation for JSON Schema. Proceedings of the VLDB Endowment, 15(13):4002-4014.
- Mohamed Amine Baazizi, Dario Colazzo, Giorgio Ghelli, Carlo Sartiani, and Stefanie Scherzinger. 2021. A json schema corpus. https://github. com/sdbs-uni-p/json-schema-corpus.
- Julian Berman. 2025. python-jsonschema. https: //github.com/python-jsonschema/jsonschema. Accessed: 2025-01-05.
- Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. 2023. Prompting Is Programming: A Query Language for Large Language Models. Proceedings of the ACM on Programming Languages, 7(PLDI):1946-1969. ArXiv:2212.06094 [cs].
- Bowtie: A meta-validator of the Bowtie. 2025. DOI: 10.5281/zenjson schema specification. odo.14646449.
- Daniel Deutsch, Shyam Upadhyay, and Dan Roth. 2019. A General-Purpose Algorithm for Constrained Sequential Inference. In Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL), pages 482-492, Hong Kong, China. Association for Computational Linguistics.
- Yixin Dong, Charlie F. Ruan, Yaxing Cai, Ruihang Lai, Ziyi Xu, Yilong Zhao, and Tianqi Chen. 2024. XGrammar: Flexible and Efficient Structured Generation Engine for Large Language Models. arXiv preprint. ArXiv:2411.15100 [cs].

Saibo Geng, Sankalp Gambhir, Chris Wendler, and Robert West. 2024. Byte bpe tokenization as an inverse string homomorphism. Preprint, arXiv:2412.03160.

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

- Saibo Geng, Martin Josifoski, Maxime Peyrard, and Robert West. 2023. Grammar-constrained decoding for structured NLP tasks without finetuning. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 10932-10952, Singapore. Association for Computational Linguistics.
- Georgi Gerganov and al. 2023. Llama.cpp: A port of facebook's llama model in c++. https://github. com/ggerganov/llama.cpp. Accessed: 2025-01-16.
- GlaiveAI. 2024. Glaive function calling dataset. https://huggingface.co/datasets/glaiveai/ glaive-function-calling. Accessed: 2024-12-21.
- Guidance AI. 2023. Guidance: A language model programming framework. Accessed: 2024-12-18.
- GuidanceAI. 2024a. Guidance acceleration tutorial. https://guidance.readthedocs.io/ en/stable/example_notebooks/tutorials/ guidance_acceleration.html. Accessed: 2025-01-16.
- GuidanceAI. 2024b. Prompt boundaries and token healing. https://github.com/guidance-ai/ guidance/blob/main/notebooks/art_of_ prompt_design/prompt_boundaries_and_ token_healing.ipynb. Accessed: 2024-12-21.
- JSON Schema Org. 2024. Json schema test suite. https://github.com/json-schema-org/ JSON-Schema-Test-Suite. Accessed: 2024-12-19.
- Kubernetes. 2022. Kubernetes json schemas. https://github.com/instrumenta/ kubernetes-json-schema. Commit hash 133f848.
- Michael Kuchnik, Virginia Smith, and George Amvrosiadis. 2023. Validating Large Language Models with ReLM. arXiv preprint. ArXiv:2211.15458 [cs].
- Will Kurt. 2024a. Coalescence: Making llm inference 5x faster. https://blog.dottxt.co/ coalescence.html. Accessed: 2024-12-21.
- Will Kurt. 2024b. Say what you mean: A response to 'let me speak freely'.
- Michael Xieyang Liu, Frederick Liu, Alexander J. Fiannaca, Terry Koo, Lucas Dixon, Michael Terry, and Carrie J. Cai. 2024. "We Need Structured Output": Towards User-centered Constraints on Large Language Model Output. In Extended Abstracts of the CHI Conference on Human Factors in Computing Systems, pages 1-9. ArXiv:2404.07362 [cs].

694

Gabriel Poesia, Oleksandr Polozov, Vu Le, Ashish Ti-

wari, Gustavo Soares, Christopher Meek, and Sumit

Gulwani. 2022. Synchromesh: Reliable code gen-

eration from pre-trained language models. arXiv

Maciej P. Polak and Dane Morgan. 2024. Extracting ac-

The Washington Post. 2022. ans-schema. https://

Subhro Roy, Sam Thomson, Tongfei Chen, Richard

github.com/washingtonpost/ans-schema. Com-

mit hash abdd6c211. Retrieved 19 September 2022.

Shin, Adam Pauls, Jason Eisner, and Benjamin Van

Durme. 2024. BenchCLAMP: A Benchmark for

Evaluating Language Models on Syntactic and Se-

mantic Parsing. arXiv preprint. ArXiv:2206.10668

Schema Store Org. 2014. The largest collection of

independent json schemas in the world. https://

www.schemastore.org/json/. A universal JSON

schema store where schemas for popular JSON docu-

ments can be found. Contributions are welcome; see

Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta

Raileanu, Maria Lomeli, Eric Hambro, Luke Zettle-

moyer, Nicola Cancedda, and Thomas Scialom. 2023.

Toolformer: Language Models Can Teach Them-

selves to Use Tools. In Advances in Neural Infor-

mation Processing Systems, volume 36, pages 68539-

Torsten Scholak, Nathan Schucher, and Dzmitry Bah-

danau. 2021. PICARD: Parsing incrementally for

constrained auto-regressive decoding from language

models. In Proceedings of the 2021 Conference on

Empirical Methods in Natural Language Processing,

pages 9895-9901, Online and Punta Cana, Domini-

can Republic. Association for Computational Lin-

Richard Shin, Christopher Lin, Sam Thomson, Charles

Chen, Subhro Roy, Emmanouil Antonios Platanios,

Adam Pauls, Dan Klein, Jason Eisner, and Benjamin Van Durme. 2021. Constrained Language Models

Yield Few-Shot Semantic Parsers. In Proceedings of

the 2021 Conference on Empirical Methods in Natu-

ral Language Processing, pages 7699-7715, Online

and Punta Cana, Dominican Republic. Association

Zhi Rui Tam, Cheng-Kuang Wu, Yi-Lin Tsai, Chieh-

Yen Lin, Hung-yi Lee, and Yun-Nung Chen. 2024.

Let me speak freely? a study on the impact of format

restrictions on large language model performance. In

Proceedings of the 2024 Conference on Empirical

Methods in Natural Language Processing: Industry

68551. Curran Associates, Inc.

for Computational Linguistics.

CONTRIBUTING.md for more information.

curate materials data from research papers with con-

versational language models and prompt engineering.

preprint. ArXiv:2201.11227 [cs].

Nature Communications, 15(1).

[cs].

701

709

- 710 711 712 713
- 714 715
- 716 717 718
- 719 720

721 722

723 724 725

726 727

729

730 731

732 733 734

735 736

737 738

739 740

741

742

743

- 744 745
- Track, pages 1218–1236, Miami, Florida, US. Asso-746 ciation for Computational Linguistics.

guistics.

Xiangru Tang, Yiming Zong, Jason Phang, Yilun Zhao, Wangchunshu Zhou, Arman Cohan, and Mark Gerstein. 2024. Struc-Bench: Are Large Language Models Good at Generating Complex Structured Tabular Data? In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers), pages 12-34, Mexico City, Mexico. Association for Computational Linguistics.

747

748

749

750

751

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

772

773

774

775

776

777

778

779

781

782

784

787

788

789

791

793

795

796

797

799

- Vivien. 2024. Llm decoding with regex constraints. https://vivien000.github.io/blog/journal/ llm-decoding-with-regex-constraints.html. Accessed: 2024-12-21.
- Bailin Wang, Zi Wang, Xuezhi Wang, Yuan Cao, Rif A Saurous, and Yoon Kim. Grammar Prompting for Domain-Specific Language Generation with Large Language Models.
- Brandon T. Willard and Rémi Louf. 2023. Efficient Guided Generation for Large Language Models. arXiv preprint. ArXiv:2307.09702 [cs].
- Austin Wright, Henry Andrews, Ben Hutton, and Greg Dennis. 2022. Draft 2020-12: Json schema core specification. https://json-schema.org/ draft/2020-12/json-schema-core.html. Published 16 June 2022. Metaschema available at https: //json-schema.org/draft/2020-12/schema.
- Shunyu Yao, Howard Chen, Austin W. Hanjie, Runzhe Yang, and Karthik Narasimhan. 2023a. COLLIE: Systematic Construction of Constrained Text Generation Tasks. arXiv preprint. ArXiv:2307.08689 [cs].
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023b. ReAct: Synergizing reasoning and acting in language models. In International Conference on Learning Representations (ICLR).
- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. 2024. SGLang: Efficient Execution of Structured Language Model Programs. arXiv preprint. ArXiv:2312.07104 [cs].

A JSON Schema Collections Details

JSONSchemaBench includes a diverse collection of schemas curated from multiple real-world applicationsAttouche et al. (2022), designed to represent a wide range of use cases:

Sources:

• GitHub (Baazizi et al., 2021): Extracted from open-source repositories containing schema definitions, representing practical, widely-used applications. Schemas from

800GitHub are of various complexities, total-
ing 6,000 schemas. We split the collec-
tion into trivial (fewer than 10 fields), easy
(10–30 fields), medium (30–100 fields), hard
(100–500 fields), and ultra (more than 500
fields), based on the total number of fields in
each JSON schema to reflect increasing com-
plexity and scale.

- **Snowplow** (**Analytics**, **2022**): Sourced from event-based analytics frameworks, showcasing schemas tailored for event-driven data structures.
 - Kubernetes (Kubernetes, 2022): Schemas defining configurations for container orchestration systems, highlighting schemas with intricate hierarchical structures.
 - WashingtonPost (Post, 2022): Schemas for The Washington Post's ANS specification.
 - GlaiveAI2K (GlaiveAI, 2024): 2,000 schemas extracted from a function-calling dataset. Each schema represents a function signature.
 - JSON Schema Store (Schema Store Org, 2014): The largest collection of independent JSON schemas in the world.

A.1 Data Processing

To ensure the quality and reliability of JSON-SchemaBench, we applied the following preprocessing steps:

1. Validation

809

810

811

812

813

814

815

816

817

818

819

821

822

823

824

826

827

830

832

836

837

838

841

- Verified schemas conform to the JSON Schema specification using the jsonschema library in Python, specifically targeting the Draft2020-12 version. Drop invalid schemas.
- Identified additional invalid schemas using validators from Rust and JavaScript libraries.

2. Cleaning

• **Deduplicate:** Removed duplicate schemas to eliminate redundancy and maintain a diverse dataset. Key ordering within schemas was ignored when determining duplicates.

• Empty Schema: Excluded schemas that were lacking meaningful constraints, effectively "empty." 844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

- Unresolved References: Removed schemas containing unresolved \$ref references to external URLs.
- Schema Version Fixes: Corrected mismatched or missing draft versions.
- Extraneous Field Removal: Eliminated unrelated fields such as command, config, path, and controls.
- **Regex Escaping:** Fixed escaping issues in regular expressions to ensure validity.
- Schema Extraction: Extracted schemas embedded within non-root levels of JSON files.

A.2 Draft versions

A.3 Feature Distribution

We count the appearance of each feature (keyword) in the 10K schemas and show the most frequent features in Figure 2a. We separately plot usage of the format keyword, which is used to specify format of string such as date-time, email, uri. This is worth highlighted because each of these formats can be quite complex to implement on its own. The distribution of formats used is shown in Figure 2b.

B Coverage Experiment Details

The prompting template used for the coverage experiment is shown in Figure 3.

Decoding Method We use greedy decoding with no top P or top K sampling for all the experiments. We only get one output from the model, which we will use to validate the schema compliance. It's totally plausible to sample more outputs and validate them all, and it might detect more schema violations. The fact that we only sample the top 1 output may quantify our *empirical coverage* as *Top 1 Empirical Coverage*.

Validation We use the jsonschema library with the Draft-2020-12 version of the JSON Schema standard to validate the generated JSON object. We turn on the 'format' checks, which are not enabled by default in Python. Strictly speaking, the jsonschema library doesn't guarantee the validation of all the schema constraints, even with

Dataset	Count	Size (KB) Med / Max	Field Count Med / Max	Max Fan-Out Med / Max	Schema Depth Med / Max
GlaiveAI-2K	1707	0.5 / 1.2	21/44	4/7	5/8
Github-Trivial	444	0.2 / 10.8	6/9	4/9	2/6
Github-Easy	1943	0.5 / 20.3	18/29	5/19	4 / 10
Snowplow	403	0.9 / 15.6	37 / 450	7 / 131	3/13
Github-Medium	1976	1.5 / 58.3	51/99	8 / 42	6/15
Kubernetes	1064	2.7 / 818.6	41/11720	5 / 600	5/7
Washington Post	125	1.7 / 81.1	44 / 2093	7 / 84	4 / 10
Github-Hard	1240	5.1 / 136.1	175 / 498	18 / 133	8 / 25
JSONSchemaStore	492	5.9 / 2934.8	155 / 108292	14 / 6543	6/22
Github-Ultra	164	25.8 / 359.6	694 / 6919	37/412	8/23

Table 8: Baisc statistics of the datasets used in the experiments.

	draft-04	draft-06	draft-07	2019-09	2020-12	unknown
Github-easy	1310	54	136	0	5	438
Github-hard	841	30	87	0	23	259
Github-medium	1221	80	140	0	7	528
JsonSchemaStore	199	5	268	5	11	4
Kubernetes	0	0	0	0	0	1087
Snowplow	0	0	0	0	0	408
WashingtonPost	125	0	0	0	0	0
Glaiveai2K	0	0	0	0	0	1707
total	4097	193	706	5	50	5155

Table 9: JSON Schema Draft Version Counts

the 'format' checks enabled. It is possible, though very rare, for a schema-noncompliant output to be validated as compliant by the jsonschema library, leading to a slight overestimation of empirical coverage. However, such occurrences are corner cases and happen infrequently.

C Theoretical Coverage Details

891

892

894

895

896

897

900

901

902

903

904

905 906

907

908

909

910

Definition C.1 (Theoretical Coverage) A

schema is considered theoretically covered if all of its features are supported by the grammar engine.

The *theoretical coverage*, noted as $C_{\text{Theoretical}}$, measures the proportion of JSON schemas that a grammar engine supports based on its implementation. It doesn't involve any model inference or experiments and is solely based on the grammar engine's implementation. $C_{\text{Theoretical}}$ is an *upper bound* of the *true coverage*, which cannot be empirically measured due to the infinite number of possible generations under the schema constraints.

Overall, the theoretical coverage provides a good indication of the grammar engine's capability to support a wide range of schema constraints.

In our experiment, the theoretical coverage for

each framework was determined based on the documentation and resources listed in Table 10. 911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

The theoretical support for each feature in JSON Schema is summarized in Figure 4

The theoretical coverage of each grammar engine is summarized in Table 11.

D JSON Schema Test Suite Experiment Details

We evaluated each constrained decoding framework's performance on the JSON Schema Test Suite using the following criteria: a framework is considered to pass a test case if it permits generating every valid instance in the test case while preventing the generation of every invalid instance. Some test cases consist exclusively of invalid instances, such as those involving *unsatisfiable* schemas, i.e., schemas for which no valid instances exist. In these cases, engines raising compile-time errors were allowed to pass.

Cleaning We removed the 'format' category of tests, as the current JSON Schema standard mandates that this keyword be ignored entirely by default. The test suite comes bundled with an 'optional' set of tests, including tests for each officially



Figure 2: Feature and Format constraint distribution.

Prompt Template for JSON Generation

System Message:

You need to generate a JSON object that matches the schema below.

Demo Examples:

941

943

947

948

951

952

953

955

```
## Input Schema: [JSON schema]
## Expected Output: [JSON object matching the schema]
...
```

Figure 3: Prompt template used to generate JSON objects in the coverage experiment.

recognized value of the 'format' keyword. We hope
to extend this work to include these optional tests
in a follow-up.

Furthermore, some tests require external resources in the form of JSON schemas available at a remote URL. We dropped these tests from the analysis, as the constrained decoding libraries discussed in the current work do not fetch these resources by default. After filtering out these tests, we are left with 43 of the original 45 test categories.

Implementation To check whether a given framework accepts or blocks the generation of a particular JSON instance, we tokenize⁵ JSONserialized form of the instance and walk the framework's constraints forward one token at a time, essentially simulating the generation process of an LLM attempting to produce the given token sequence:

> • XGrammar directly expose an interface for updating the token mask after inserting a token and checking validity.

• Outlines does not expose a public interface for interacting with the token mask, but outlines-core, which outlines is built on top of, is easily adapted for this purpose. 956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

- Similarly, Guidance does not expose a public interface for interacting with the token mask, but llguidance, which guidance is built on top of, is easily adapted for this purpose.
- Llamacpp does not expose this interface, but it shares a common grammar-specification language with XGrammar. We use llamacpp to generate GGML BNF and check tokensequence validity using xgrammar's interface.

We provide code snippets that show the use of the JSON Schema Test Suite to assess the test coverage of each constrained decoding framework. For each framework, we implemented a 'test harness' according to the base classes showed in listing 1.

Listing 2 shows the criteria for a test case to pass, which depends on all tests in the case to pass (listing 3). We show the definition of TestCase and Test in listing 4.

Concrete implementations of the test harness for each framework are reported in listings 5, 6, 7, and

⁵The particular choice of tokenizer is not particularly important, but we use the Llama 3.1 tokenizer for consistency with our other experiments.

Lib Version	Release Date	JSON Schema Support Documentation
0.2.0rc	2024.11.26	LLGuidance Documentation
0.3.2	2024.11.16	llama.cpp JSON Schema to gbnf Conversion
0.1.6	2024.12.07	XGrammar JSON Schema to gbnf Conversion
0.1.8	2024.12.06	Outlines JSON Schema to Regex Conversion
UNK	UNK	OpenAI Structured Output API
0.8.3	2024.10.31	Gemini Structured Output Content Types
	Lib Version 0.2.0rc 0.3.2 0.1.6 0.1.8 UNK 0.8.3	Lib VersionRelease Date0.2.0rc2024.11.260.3.22024.11.160.1.62024.12.070.1.82024.12.06UNKUNK0.8.32024.10.31

Table 10: Grammar Engine Documentation and Resources

Dataset	LM only	Guidance	Llamacpp	Outlines	XGrammar	OpenAI	Gemini
GlaiveAI	0.00	0.96	0.95	0.95	0.87	0.87	0.87
GitHub Easy	0.00	0.87	0.83	0.75	0.65	0.31	0.31
Snowplow	0.00	0.80	0.74	0.58	NA	0.29	NA
GitHub Medium	0.00	0.73	0.69	0.57	0.49	0.22	NA
Kubernetes	0.00	0.58	0.58	0.58	0.58	0.40	NA
Washington Post	0.00	0.70	0.64	0.63	0.62	0.29	NA
GitHub Hard	0.00	0.54	0.49	0.38	0.33	0.00	NA
JsonSchemaStore	0.00	0.31	0.24	0.20	0.13	0.00	NA

Table 11: Theoretical coverage across datasets.

981

982

991

993

995

997

998

999

1001

1002

1003

1004 1005

1006

1007

1009

8.

E Efficiency Experiment Details

For efficiency experiments, the results depend on both the size of the model and the tokenizer's vocabulary size. We used **Llama-3.1-8B-Instruct** (quantized to Q8bit) with a 128K token vocabulary to achieve a balance between computational efficiency and model capability.

Below, we outline specific considerations related to grammar and prefix caching:

 Grammar Cache (Compilation): Since each schema in the dataset is unique, caching grammar compilations does not offer any benefits.

• **Prefix Cache (LLM Inference):** We implement prefix caching during LLM inference for all cases to enhance efficiency by reusing computed results where applicable.

F Quality Experiment Details

Prompt and JSON Schema For the task of **Shuf-fle Objects**, and **GSM8K**, we use the same prompt and JSON schema from the dottxt's "let me speak freely" rebuttal.

For the task of **Last Letter**, we make a slight modification because the original prompt used was a bad example as pointed out by (Kurt, 2024b). We also put it into a JSON format to better align with the other tasks.

Figure 8 reveals non-empty exclusive regions for each engine, indicating that no single engine outperforms the others across all instances.

G Engine calling Snippet

We provide a snippet of the engine code used1011in our experiments. The generation method1012of each engine has two main components:1013"compile_grammar" and "call_engine".1014

1010

1015

1016

1017

1018

1019

1020

1021

We use the Listing 10 to validate the generated JSONs against the schema. The validation is done by the jsonschema library with format checking enabled.

We provide a snippet of how the engines are called in our experiments in Listings 11, 12, 13, and 14.

Feature	Schemas	LLGuidance	llama.cpp	Outlines	XGrammar	OpenAl	Gemini
required	7473					×	
items	4703						
additionalProperties	4290					×	×
enum	3898						
\$ref	2786						
pattern	1829					×	×
format	1570				×	×	\checkmark
oneOf	1437	×	×	×	×	×	×
@minmaxLength	1388				×	×	×
@siblingKeys	1230		×	×	×	×	×
@minmaxInteger	1217			×		×	×
@minmaxItems	1132				×	×	×
additionalProperties:object	1086					×	×
anyOf	924						
patternProperties	783	×	×	×	×	×	×
allOf	756			×	×	×	×
not	670	×	×	×	×	×	×
@minmaxProperties	235	×	×		×	×	×
additionalItems	221					×	×
const	183						×
dependencies	156	×	×	×	×	×	×
multipleOf	134	×	×	×	×	×	×
@minmaxNumber	123		×	×	×	×	×
uniqueltems	117	×	×	×	×	×	×
if	63	×	×	×	×	×	×
propertyNames	56	×	×	×	×	×	×
contains	23	X	×	×	×	×	×
unevaluatedProperties	2	×	×	×		×	×
@recursiveSchemas	0			×			×
Schemas supported	10206	7100	6727	6081	5401	3403	3401

Figure 4: Feature checklist for different structured output engines



Figure 5: Illustration of over-constrained and under-constrained.

Dataset	Framework	GCT (s)	TTFT (s)	TPOT (ms)	TGT (s)	Output Tokens (FF)
GlaiveAI	LLM only	NA	<u>0.10</u>	15.40	1.08	64.94 (00.00)
	Guidance	<u>0.00</u>	0.24	<u>6.37</u>	<u>0.50</u>	41.56 (15.70)
	Llamacpp	0.05	0.20	29.98	1.47	43.18 (00.00)
	Outlines	3.48	3.65	30.33	4.84	40.39 (00.00)
GitHub Easy	LLM only	NA	<u>0.10</u>	15.83	0.95	53.91 (00.00)
	Guidance	<u>0.00</u>	0.34	<u>7.44</u>	<u>0.60</u>	34.92 (10.02)
	Llamacpp	0.05	0.18	27.22	1.10	33.93 (00.00)
	Outlines	3.71	3.97	39.78	5.29	34.19 (00.00)
Snowplow	LLM only	NA	<u>0.11</u>	16.23	1.01	55.31 (00.00)
	Guidance	<u>0.00</u>	0.28	<u>6.55</u>	<u>0.51</u>	36.77 (14.50)
	Llamacpp	0.05	0.20	28.90	1.24	37.21 (00.00)
	Outlines	3.91	4.14	42.66	5.65	35.65 (00.00)
GitHub Medium	LLM only	NA	0.20	16.68	2.56	142.10 (00.00)
	Guidance	<u>0.01</u>	0.54	<u>7.57</u>	<u>1.29</u>	99.66 (31.42)
	Llamacpp	0.06	0.30	29.08	2.85	87.71 (00.00)
	Outlines	8.05	8.38	46.57	12.23	84.64 (00.00)
Kubernetes	LLM only	NA	<u>0.16</u>	15.32	0.84	44.38 (00.00)
	Guidance	<u>0.01</u>	0.45	<u>9.47</u>	<u>0.71</u>	28.75 (04.40)
	Llamacpp	0.05	0.28	28.04	1.06	28.09 (00.00)
	Outlines	5.29	5.55	46.10	6.56	22.26 (00.00)

Table 12: **Efficiency metrics** for different engines with **LlamaCpp** as the inference engine. **GCT**: Grammar Compilation Time, **TTFT**: Time to First Token, **TPOT**: Time Per Output Token, **TGT**: Total Generation Time, **FF**: Fast-Forwarded output tokens. Bold values indicate the smallest in each column for GCT, TTFT, TPOT, and TGT. All values are **median** of the samples.



JSON Schema Test Suite Coverage Framework

Figure 6: JSON Schema test suite coverage by category. Each cell represents the proportion of passing tests for a given category-framework pair, with darker shades indicating higher coverage. A single asterisk (*) marks frameworks tied for the highest (non-zero) coverage, while a double asterisk (**) marks the framework with the single highest coverage in the category.

```
class Compiler:
    def __init__(self, model_id: str):
        .....
        Builds a Compiler, taking a huggingface model_id to provide
        configuration information about the model and/or tokenizer.
        .....
    def compile(self, schema: str) -> Masker:
        .....
        Compiles a schema into a masker used to validate a stream of
        tokens according to the schema.
        Raises an exception if the framework cannot compile the schema.
        .....
class Masker:
    def advance(self, token: int):
        .....
        Advances the masker by one token.
        Raises an exception if the token is not allowed by the mask.
        .....
    def assert_done(self):
        .....
        Asserts that the masker is either in a terminal state or will
        accept an EOS token, after which it will be in a terminal state.
```

Raises an exception if otherwise.

Listing 1: Abstract test harness

Dataset	Framework	GCT (s)	TTFT (s)	TPOT (ms)	TGT (s)	Output Tokens (FF)
GlaiveAI	Guidance	<u>0.01</u>	0.36	<u>36.92</u>	<u>1.87</u>	41.45(16.76)
	XGrammar	0.12	<u>0.30</u>	66.78	2.87	39.47(00.00)
GitHub Easy	Guidance	<u>0.01</u>	0.37	<u>42.03</u>	<u>1.60</u>	27.67(06.75)
	XGrammar	0.11	<u>0.33</u>	65.57	4.07	59.45(00.00)
GitHub Medium	Guidance	<u>0.01</u>	0.55	<u>44.21</u>	<u>4.84</u>	96.31(26.93)
	XGrammar	0.20	<u>0.48</u>	65.51	6.53	92.93(00.00)
GitHub Hard	Guidance	<u>0.01</u>	0.73	<u>35.88</u>	<u>10.25</u>	211.40(101.40)
	XGrammar	0.30	0.65	65.20	14.99	221.40(00.00)

Table 13: Efficiency metrics for different engines with Hugging Face Transformers as the inference engine. All values are median of the samples.

```
def do_test_case(test_case: TestCase, compiler: Compiler, tokenizer: Tokenizer) ->
    bool:
\hookrightarrow
    try:
        masker = compiler.compile(json.dumps(test_case.schema))
    except:
        if all(not test.valid for test in test_case.tests):
            # Pass: compile error on a case with only invalid test data
            return True
        else:
            # Fail: compile error but schema has at least one valid test datum
            return False
    for test in test_case.tests:
        passed = do_test(test, tokenizer, masker.copy())
        if not passed:
            # Fail: a test failed
            return False
    # Pass: all tests passed
    return True
```

Listing 2: Running a test case

```
def do_test(test: Test, tokenizer: Tokenizer, masker: Masker) -> bool:
    tokens = tokenizer(json.dumps(test.data),
    → add_special_tokens=False)["input_ids"]
    try:
        for token in tokens:
            masker.advance(token)
        masker.assert_done()
    except:
        if test.valid:
            # Fail: valid data was rejected
            return False
        else:
            # Pass: invalid data was rejected
            return True
    else:
        if test.valid:
            # Pass: valid data was accepted
            return True
        else:
            # Fail: invalid data was accepted
            return False
```

Listing 3: Running a test

```
from pydantic import BaseModel
from typing import Any, Union
class TestCase(BaseModel):
    schema: Union[bool, dict]
    tests: list[Test]
class Test(BaseModel):
    data: Any
    valid: bool
                              Listing 4: TestCase specification
import outlines
import outlines_core
class OutlinesCompiler(Compiler):
    def __init__(self, model_id: str):
        self.tokenizer = outlines.models.transformers(model_id).tokenizer
    def compile(self, schema: str) -> "OutlinesMasker":
        regex = build_regex_from_schema(schema)
        guide = outlines.fsm.guide.RegexGuide.from_regex(
            regex, self.tokenizer
        )
        return OutlinesMasker(guide, eos_token_id=self.tokenizer.eos_token_id)
class OutlinesMasker(Masker):
    def __init__(self, guide, eos_token_id=None):
        self.guide = guide
        self.state = self.guide.initial_state
        self.eos_token_id = eos_token_id
    def advance(self, token: int):
        assert token in self.guide.get_next_instruction(self.state).tokens
        self.state = self.guide.get_next_state(self.state, token)
    def assert_done(self):
        if not self.guide.is_final_state(self.state):
            assert self.eos_token_id in
            → self.guide.get_next_instruction(self.state).tokens
            self.advance(self.eos_token_id)
        assert self.guide.is_final_state(self.state)
```

Listing 5: Concrete test harness for Outlines

```
import guidance
import llguidance
class GuidanceCompiler(Compiler):
    def __init__(self, model_id: str):
        self.gtokenizer =
        → guidance.models.transformers.TransformersTokenizer(model_id, None)
        self.lltokenizer =
        → llguidance.LLTokenizer(llguidance.TokenizerWrapper(self.gtokenizer))
    def compile(self, schema: str) -> GuidanceMasker:
        grammar = guidance.json(schema=schema)
        llinterpreter = llguidance.LLInterpreter(
            tokenizer=self.lltokenizer,
            llguidance_json=json.dumps(grammar.ll_serialize()),
            enable_backtrack=False,
            enable_ff_tokens=False,
        )
        return GuidanceMasker(llinterpreter, self.gtokenizer.eos_token_id)
class GuidanceMasker(Masker):
    def __init__(self, llinterpreter, eos_token_id):
        self.llinterpreter = llinterpreter
        self.eos_token_id = eos_token_id
    def advance(self, token: int):
        bytemask, _ = self.llinterpreter.compute_mask()
        assert bytemask[token] > 0
        self.llinterpreter.commit_token(token)
    def assert_done(self):
        if self.llinterpreter.stop_reason() == "NotStopped":
            bytemask, _ = self.llinterpreter.compute_mask()
            if bytemask is not None:
                assert bytemask[self.eos_token_id] > 0
                self.llinterpreter.commit_token(self.eos_token_id)
                bytemask, _ = self.llinterpreter.compute_mask()
                assert bytemask is None
        assert self.llinterpreter.stop_reason() in {"NoExtension", "EndOfSentence"}
```

Task	Example	Structure	Metric
Last Letter	Input: Ian Peter Bernard Stephen Output: nrdn	CoT reasoning + answer in $a - z$	Case-sensitive ex- act match
Shuffle Objects	Input: Sequence of exchanges among indi- viduals + choices Output: A-E	CoT reasoning + answer in $A - E$	Exact match
GSM8K	Input: Basic calculation problems Output: Number, e.g., 8	CoT reasoning + an- swer as integer	Number exact match

Listing 6: Concrete test harness for Guidance

Table 14: Task Descriptions and Structures

```
import xgrammar as xgr
from transformers import AutoConfig, AutoTokenizer
class XGrammarCompiler(Compiler):
    def __init__(self, model_id: str):
        tokenizer = AutoTokenizer.from_pretrained(model_id)
        config = AutoConfig.from_pretrained(model_id)
        self.eos_token_id = tokenizer.eos_token_id
        self.tokenizer_info = xgr.TokenizerInfo.from_huggingface(
            tokenizer, vocab_size=config.vocab_size
        )
        self.compiler = xgr.GrammarCompiler(
            tokenizer_info=self.tokenizer_info,
        )
    def compile(self, schema: str) -> "XGrammarMasker":
        compiled_grammar = self.compiler.compile_json_schema(schema,
        → strict_mode=False)
        xgr_matcher = xgr.GrammarMatcher(compiled_grammar)
        return XGrammarMasker(xgr_matcher, self.eos_token_id)
class XGrammarMasker(Masker):
    def __init__(self, xgr_matcher, eos_token_id):
        self.matcher = xgr_matcher
        self.eos_token_id = eos_token_id
    def advance(self, token: int):
        assert self.matcher.accept_token(token)
    def assert_done(self):
        if not self.matcher.is_terminated():
            self.advance(self.eos_token_id)
        assert self.matcher.is_terminated()
                         Listing 7: Concrete test harness for xGrammar
from llama_cpp import LlamaGrammar
```

```
import xgrammar as xgr
```

```
class LlamacppCompiler(XGrammarCompiler):
```

```
def compile(self, schema) -> XGrammarMasker:
    grammar_bnf = LlamaGrammar.from_json_schema(schema)._grammar
    compiled_grammar = self.compiler.compile_grammar(grammar_bnf)
    xgr_matcher = xgr.GrammarMatcher(compiled_grammar)
    return XGrammarMasker(xgr_matcher, self.eos_token_id)
```

Listing 8: Concrete test harness for Llamacpp, inheriting from the XGrammar harness for all functionality after using llamacpp to convert the schema to GGML BNF.

```
import time
import stopit
class BaseModel:
   @stopit.threading_timeoutable(timeout=40)
   def compile_grammar(self, json_schema):
        status = "unknown"
        try:
            compiled_grammar = self._compile_grammar(json_schema)
            status = "success"
        except Exception as e:
            # Any exception in this block will be caught and considered as schema
            \rightarrow not supported
            compiled_grammar = None
            status = "schema_not_supported"
        return compiled_grammar, status
   def generate(self, prompt, json_schema=None):
        compile_start_time = time.time()
        compiled_grammar = self.compile_grammar(json_schema)
        compile_end_time = time.time()
        # GCT (Grammar Compilation Time)
        gct = compile_end_time - compile_start_time
        gen_start_time = time.time()
        output, first_tok_arr_time = self._call_engine(prompt, compiled_grammar)
        # TTFT (Time to First Token)
        ttft = first_tok_arr_time - gen_start_time
        gen_end_time = time.time()
        # TGT (Total Generation Time)
        tgt = gen_end_time - gen_start_time
        return output, gct, ttft, tgt
   def _call_engine(self, prompt, compiled_grammar):
        raise NotImplementedError
```

Listing 9: Abstract BaseModel interface defining the calling of structured generation, including grammar compilation and text generation timing metrics.

```
import jsonschema
```

```
from jsonschema import Draft202012Validator, FormatChecker, ValidationError
format_checker = FormatChecker()

def is_json_schema_valid(schema: dict):
    try:
        jsonschema.Draft202012Validator.check_schema(schema)
        return True
    except jsonschema.SchemaError as e:
        return False

def validate_json_against_schema(json_obj, json_schema):
    if not is_json_schema_valid(json_schema):
        raise ValidationError("The JSON schema is invalid.")
        validator = Draft202012Validator(json_schema, format_checker=format_checker)
        return validator.validate(json_obj)
```

Listing 10: Validation of the generated JSONs against the schema.

import guidance

```
class GuidanceModel(BaseModel):
    def compile_grammar(self, json_schema):
        return guidance.json(
            schema=json_schema,
        )
    def _call_engine(self, prompt, compiled_grammar):
        generator = self.guidance_model.stream() + prompt + compiled_grammar
        for i, state in enumerate(generator):
            if i == 0:
                first_state_arr_time = time.time()
        output = state
        return output, first_state_arr_time
```

Listing 11: Invocation of the guidance engine.

import llama_cpp

Listing 12: Invocation of the LlamaCpp engine.

import outlines

Listing 13: Invocation of the Outlines engine.

import xgrammar

```
class TimingLogitsProcessor(LogitsProcessor):
    def __init__(self):
        super().__init__()
        self.timestamps = []
    def __call__(self, input_ids, scores):
        current_time = time.time()
        self.timestamps.append(current_time)
        return scores
class XGrammarModel(BaseModel):
    def compile_grammar(self, json_schema):
        return xgrammar.GrammarCompiler().compile_json_schema(json_schema)
    def _call_engine(self, prompt, compiled_grammar):
        output = self.hf_model.generate(prompt, logits_processor=[compiled_grammar,

    timeit_logit_processor])

        first_tok_arr_time = timeit_logit_processor.timestamps[0]
        return output, first_tok_arr_time
```

Listing 14: Invocation of the XGrammar engine.

Prompt Template for GSM8K

System Message:

You are an expert in solving grade school math tasks. You will be presented with a grade-school math word problem and be asked to solve it. Before answering, you should reason about the problem (using the "reasoning" field in the JSON response format described below). Always respond with JSON in the format: {"reasoning": <reasoning about the answer>, "answer": <final answer>}. The "reasoning" field contains your logical explanation, and the "answer" field contains the final numeric result.

Demo Examples:

```
## Input: "[example question]"
## Output: "reasoning": "[example reasoning]", "answer": [example answer]
```

•••

Figure 7: Prompt template for solving GSM8K with JSON responses.



Figure 8: Overlap of Correct Instances Across Models on GSM8K