# FANTAstic SEquences and Where to Find Them: Faithful and Efficient API Call Generation through State-tracked Constrained Decoding and Reranking

Anonymous ACL submission

#### Abstract

001 API call generation is the cornerstone of large language models' tool-using ability that provides access to the larger world. However, existing supervised and in-context learning approaches suffer from high training costs, poor data efficiency, and generated API calls that can be unfaithful to the API documentation and the 800 user's request. To address these limitations, we propose an output-side optimization approach called FANTASE. Two of the unique contributions of FANTASE are its State-Tracked Con-011 strained Decoding (SCD) and Reranking components. SCD dynamically incorporates appropriate API constraints in the form of Token Search Trie for efficient and guaranteed generation faithfulness with respect to the API documentation. The Reranking component effi-017 ciently brings in the supervised signal by leveraging a lightweight model as the discriminator to rerank the beam-searched candidate generations of the large language model. We demon-021 strate the superior performance of FANTASE in API call generation accuracy, inference efficiency, and context efficiency with DSTC8 and API Bank datasets.

## 1 Introduction

026

027

028

In recent year, there has been a surge of interest in enabling the automated tool-using capability of intelligent systems (Schick et al., 2023; Mialon et al., 2023). Specifically, as a bridge to the larger world, Application Programming Interface (API) calls allow virtual assistants to control smart-home devices, retrieve information, make reservations, and more on the user's behalf. Figure 1 shows how an API call may improve the user-assistant conversation and satisfy the user's needs. Generating such an API call requires advanced capabilities in understanding the requirements of an API (including its endpoints, parameters, and expected data formats) and reasoning over the conversation context to translate the user's needs into the appropriate API format.



Restaurants\_1.FindRestaurants(city="Mountain View", cuisine="American", price\_range="moderate")

Figure 1: Example of an API call that retrieves information based on the user's needs given in the conversation.

047

054

055

060

061

062

063

064

065

066

067

068

069

With recent breakthroughs in generative Large Language Models (LLMs) such as GPT-X (Ouyang et al., 2022; OpenAI, 2023) and LLaMA (Touvron et al., 2023a,b), researchers have started to investigate their competence in complex reasoning tasks such as utilizing appropriate API tools (Li et al., 2023; Qin et al., 2023; Wang et al., 2023a). Their attempts focus on methods that can generally be grouped into those based on *supervised fine-tuning* for task-specific usage and those based on augmenting input-side context information (such as API shortlisting and exemplar selection) and optimizing prompts for in-context learning (Brown et al., 2020; Wei et al., 2022). Despite the strong supervision or extensive context, these methods still cannot ensure the generation's faithfulness with respect to the API documentation and suffer data and compute inefficiency. In contrast to previous works, we focus on how decoding strategies improve the generation's faithfulness, which is complementary to supervised fine-tuning and in-context learning methods. As a result, we present FANTASE (FAN-TAstic SEquences and Where to Find Them), a framework that employs State-tracked Constrained Decoding (SCD) and Reranking components, for faithful and efficient API call generation.

The SCD tracks the states of the generation and retrieves appropriate API documentation con-

straints in the form of Constrained Token Search 071 Trie (CTST) used at each decoding step. SCD 072 is guaranteed to generate API calls that are faith-073 ful with respect to the API documentation ( $\S$  4.1), and provides inference efficiency (§ 6.2) with CTST that eliminates unnecessary forward inference passes. Compared to supervised fine-tuning 077 methods, SCD brings considerable improvements  $(\S 6.1)$  without the data labeling and model training related hefty costs of labor, time, and computing that become increasingly expensive as the size of the LLM grows (Yang et al., 2023). SCD also reduces the in-context learning's reliance on the repeated supply of extensive contextual information 084 for the inference of each instance  $(\S 6.3)$  by effectively incorporating API documentation constraints and guaranteeing the associated faithfulness at the decoding stage.

The Reranking component of FANTASE leverages models that are significantly smaller than 091 LLMs for efficient incorporation of supervised signals ( $\S$  4.2). As the correct API generation may not always have the highest sequence probability among beam-searched candidate sequences (§ 3), we train lightweight models to discriminate and rerank LLMs' candidate generations and demonstrate their effectiveness in digging out those cor-097 rect sequences ( $\S$  6.1). Compared to the supervised fine-tuning of LLMs, the Reranking component features extremely low training costs as it employs 100 lightweight models. Compared to input-side op-101 timized in-context learning methods, the Rerank-102 103 ing component can address the severe performance issue associated with the absence of valuable su-104 pervised signals. Notably, FANTASE is a highly 105 adaptable approach that suits the evolving and vast nature of real-world APIs. With the update of API 107 documentation or the application to the new do-108 main, LLMs fine-tuned with old data would require 109 re-tuning with new data (Kumar et al., 2022). For 110 FANTASE, SCD can easily adapt by constraining 111 the decoding with a new set of constraints elicited 112 from the new API documentation, while re-tuning 113 the lightweight Reranking models has lower time 114 and compute cost. 115

In summary, we make the following novel contributions:

116

117

118

119

120

121

• We propose State-tracked Constrained Decoding that can effectively enforce constraints elicited from API Documentation, which yields faithful generation and context efficiency.

 We leverage Constrained Token Search Trie to reduce unnecessary forward inference passes, which yields faster generation speed.
 We demonstrate the effectiveness of incorporating supervised signals with a small model

122

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

rating supervised signals with a small model to discriminate and rerank the beam-searched candidate generations of LLMs.

## 2 Related Work

Constrained Decoding offers controllable text generation by enforcing certain constraints at the decoding stage. Early research (Hokamp and Liu, 2017; Post and Vilar, 2018) concentrated on lexical constraints that enforce the inclusion of specific words or phrases in the outputs, which often neglects broader syntactic or semantic relationships. Later on, Lu et al. 2021 introduced *NeuroLogic* Decoding that handles more complex lexical constraints expressed by predicate logic. The subsequent extension, NeuroLogic A\*esque Decoding (Lu et al., 2022), incorporated a lookahead heuristic to estimate future lexical constraint satisfaction. More recently, Chen et al. 2022 and Bastan et al. 2023 proposed parsing-based constrained decoding algorithms that tackle the challenge of ensuring correct syntactic relationships between word pairs.

Specific to structured text generation, Scholak et al. 2021 targeted Text-to-SQL generation and introduced PICARD that checks the validity at each decoding step for SQL lexical and grammar correctness with incremental parsing. The latest advancement was made by Geng et al. 2023 who demonstrated that an incremental parser can be used with formal grammar on a much wider range of structured NLP tasks without finetuning. While the results are encouraging, existing methods require post-hoc constraint satisfaction checking or rely on dependency parsing at inference time, or both, which compromises the efficiency. The most recent and closest work to ours is API-aware Constrained Decoding (Wang et al., 2023a) that imposes function and argument token constraints based on API documentation. However, despite limited improvements, its decoding strategy results in a 20% slowdown of the generation. In contrast to aforementioned methods, we achieve faster generation speed and guaranteed faithfulness with a novel Statetracked Constrained Decoding approach that dynamically incorporates appropriate constraints in

Deleted	Human: I want to find a burger joint. Assistant: In which city?.
Conversation	Human: In Mountain View. Assistant: Eureka! restaurant is in Mountain View.
	Human: Are there any other restaurants in the moderate price range?
Delated	Restaurants_1.FindRestaurants("cuisine" : Required, "city" : Required, "price_range" : Optional,
A PI Documentation	"has_live_music" : Optional, "serves_alcohol" : Optional) the possible values for "cuisine" include
Ar i Documentation	["Mexican", "Chinese", "Indian", "American", "Italian"]
Expected API Call	Restaurants_1.FindRestaurants(city="Mountain View", cuisine="American", price_range="moderate")
	1. Restaurants_1.FindRestaurants(price_range="moderate", city="MountainView") missing cuisine
Top Candidates	2. Restaurants_1.FindRestaurants(cuisine="Burgers", city="MountainView") missing price_range
	3. Restaurants_1.FindRestaurants(cuisine="American", city="MountainView") missing price_range
	<ol><li>Restaurants_1.FindRestaurants(price_range="moderate") missing cuisine and city</li></ol>
	<ol><li>Restaurants_1.FindRestaurants(cuisine="American", city="MountainView", price_range="moderate")</li></ol>

Table 1: Preliminary analysis sample. With regular beam search decoding, the correct generation is only ranked the 5th, and other higher ranked generations exhibit various errors highlighted in red.

the form of a retrieved token search trie.

172

204

205

208

Discriminator Guided Generation utilizes small 173 discriminative models or external tools to guide the 174 generation of LLMs. Dathathri et al. 2020 proposed 175 the Plug and Play Language Model concept that 176 guides the generation of pretrained models with 177 a lightweight attribute classifiers' gradient. How-178 ever, it increases compute costs due to the extra for-179 ward and backward passes required for sampling 180 and using the gradients from the attribute classi-181 fiers to push the pretrained model's hidden activations. Following works including GeDi (Krause et al., 2021), FUDGE (Yang and Klein, 2021), 184 and BeamR (Landsman et al., 2022) used different lightweight discriminators that classify the attribute of possible next tokens or partial sequence 187 and reweigh token-level or beam-level probabilities at each decoding step towards the desired direction 189 190 of attributes like sentiment, topic, formality, and so on. More recently, Ni et al. 2023 leveraged the 191 execution results of a SQL executor to steer SQL 192 generation, which achieved new state-of-the-art results. Nevertheless, the method is bounded by the 194 195 prerequisite of the external executor. In our work, we employ a lightweight model to discriminate API 196 call generation by the given context and perform 197 a one-pass reranking of the beam-searched results, 198 which brings in supervised signals effectively with 199 little compute and time costs to the overall generation framework. 201

## **3** Preliminary Analysis

To better understand the capabilities and limitations of existing LLMs on the task of API call generation, we conduct a preliminary inference analysis on one hundred DSTC8 (Kim et al., 2019)<sup>1</sup> samples with an Alpaca (Taori et al., 2023) model that had been tuned with GPT-generated self-instruct (Wang et al., 2023b) data for better instruction following and in-context learning capabilities. We prompt the model with the DSTC8 data that contains task instruction, documentation of related APIs, two related exemplars, and conversation history. We use beam search with beam size 10 as the decoding algorithm, and we consider the top-10 high probability sequences as the candidate generations.

209

210

211

212

213

214

215

216

217

218

219

221

222

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

240

241

242

243

244

245

246

247

Our quantitative analysis shows that for 73% of the cases, the correct API calls are generated within those high probability sequences. However, within these cases, almost half of the correct sequences were not ranked as the highest, which yields a top-1 API call generation accuracy of 41%. Table 1 presents an example where the user wants to find a burger joint with a moderate price range in Mountain View. The supplied API documentation specified that the Restaurants\_1.FindRestaurants function has cuisine and city as the required arguments, and the cuisine argument has five possible values. However, the correct sequence was only ranked the 5th for the given example. All the other 4 candidates that have higher sequence probabilities missed some required arguments, and the second one also wrongly generated Burgers instead of one of the five possible values for the argument cuisine. Note that the model demonstrates some reasoning capability that can correctly map Burgers into American as shown in the second and the fifth sequences. Nevertheless, the overall sequence probability favors the problematic generation of Burgers, which may be attributed to the explicit mention of the word in the given conversation history.

We conduct a further qualitative analysis to categorize the error types and possible mitigation for these cases. For the highest-ranked error cases, we find 33% argument value error, 24% missing required arguments, 19% missing optional arguments, 14% hallucination, and 10% argument name

<sup>&</sup>lt;sup>1</sup>Details will be given in Section 5.1



Figure 2: Illustration of the Concepts of Constrained Decoding and Reranking. (Upper half) Constrained Decoding enforces API documentation constraints and would only consider the five possible values of cuisine. (Lower half) A lightweight RoBERTa model is used to discriminate and rerank the beam searched candidate generations.

error. Furthermore, 42% of the error cases can be mitigated by enforcing the constraints described in the API documentation, 29% of the cases require the better understanding of the conversation, and the remaining 29% of the cases need a combination of the aforementioned two improvements.

### 4 The FANTASE Framework

249

254

258

259

260

264

269

271

272

In Section 3, we demonstrate that there are "FAN-TAstic SEquences" in the beam-searched candidate generations, and the question is where to find them. To dig out those "FANTAstic SEquences" with the data efficiency and compute efficiency in mind, we propose the FANTASE framework that consists of two major compotents – State-tracked Constrained Decoding (§4.1) and Reranking (§4.2), which aims at enforcing API constraints with guaranteed faithfulness to the API documentation and incorporating supervised signals at low compute costs respectively.

#### 4.1 State-tracked Constrained Decoding

In Figure 2, we illustrate the concept of Statetracked Constrained Decoding (SCD). For a regular decoding step, the consideration of the entire vocabulary space would lead to the high probability of the word Burgers overshadowing the correct word of American. To ensure the faithfulness to the API documentation, our SCD approach enforces the model to only consider the probabilities of the five possible values as documented in the API documentation.

Different from conventional token-occurrencebased constrained decoding approaches as described in Section 2, our approach takes the relation between package, function, argument, and argument values into consideration. SCD allows precise and dynamic enforcement of constraints based on the API documentation and generated units, which avoids the look-ahead decoding and pruning as other constrained decoding algorithms would normally require. The implementation of SCD consists of three major parts: 1) extraction of constraints from API documentation, 2) state tracking of the generation for constraints retrieval, and 3) constrained decoding with token search trie.

As API documentation is usually well-structured, it is feasible to extract constraints with simple rules. As a preprocessing step, we use regular expressions to extract five types of constraints including 1) available packages, 2) functions of each package, 3) required arguments of each function, 4) optional 273

274



Restaurants\_1.FindRestaurants(cuisine="Burgers", city="Mountain View") missing price\_range

Figure 3: State-tracked Constrained Generation of API Call (showing the step of generating the value of parameter cuisine that has possible values of American, Chinese, Indian, Italian, and Mexican).

Structural Token	Generation State	Actions
S (poudo)	Start of the conception	- Retrieve all package names
is (psudo)	Start of the generation	- Constrained generation of package name
	End of peakage name	- Record decoded package name
DOT	Stort of function name	- Retrieve possible function names by using package name
DOI	Start of function name	- Constrained generation of function name
(	End of function name	- Record decoded function name
LEET BRACKET	Start of argument nome	- Retrieve possible argument names by using package and function names
LEFT_BRACKET	Start of argument name	- Constrained generation of argument name
		- Record decoded argument name
_	End of argument name	- Retrieve possible argument values by using package, function, and argument names
FOUAL	Stort of argument value	- Check if the argument only takes certain possible values
LQUAL	Start of argument value	- If so, constrained generation of argument value
		- If not, perform normal unconstrained generation
,	End of argument value	- Reuse previously retrieved possible argument names
COMMA	Start of argument name	- Constrained generation of argument name
)		- Check if the list of decoded argument name contains all the required arguments
RIGHT BRACKET	End of the generation	<ul> <li>If so, conclude the generation</li> </ul>
KIGHT_BRACKET		- If not, replace RIGHT_BRACKET with COMMA and enforce continued generation

Table 2: State Tracking with Structural Tokens and Associated Actions.

arguments of each function, and 5) possible values 299 of each argument. We store these constraints in a lookup table with package name, function name, 301 and argument name as the query keys. At the inference stage, we query the lookup table to fetch cor-303 responding constraints by decoded package name, 304 function name, and/or argument name. If the API is 305 evolved with changing constraints such as new re-307 quired/optional arguments, changing names, changing possible values, etc., new constraints can be enforced effortlessly at the inference stage by reparsing the updated API documentation, which is less expensive than re-tuning the model with up-311 dated labeled data. 312

In Figure 3, we illustrate the SCD at the infer-313 ence stage. To ensure appropriate constraints can be retrieved and enforced at the precise inference 315 step of the generation, the state of the generation is 316 determined by tracking the model-generated struc-317 tural tokens. The structured nature of the API call results in signature tokens that indicate the end or 319

start of different units of the API call as we speci-320 fied in Table 2. Specifically, for the constrained generation of an API call unit, we enforce the model 322 to decode along the Constrained Token Search Trie (CTST) as illustrated in Figure 4. The tokenizer 324 of LLMs performs WordPiece tokenization, which 325 breaks down the word into smaller subword tokens for various benefits (Devlin et al., 2019). Accordingly, the package name Restaurants\_1 would 328 be autoregressively generated by the LLM piece by piece with five forward-pass inference steps as 330 shown in the upper half of Figure 4. At the prepro-331 cessing step, we build the extracted constraints into 332 CTST. When conducting constrained generation, the forward inference pass is only necessary for 334 nodes that have multiple branches. In such a case, only the probabilities of the possible next tokens as 336 indicated by the CTST would be considered. For 337 nodes that only have one child, the subsequent token is directly appended, which saves the time and compute costs of a forward inference pass.

321

326

371

341

343



Figure 4: Comparison of the Generation of Restaurants\_1 with Normal Decoding and Decoding with Constrained Token Search Trie.

For SCD, we implement it with four different sampling strategies including greedy search, top-k sampling, top-p sampling, and beam search.

#### 4.2 Reranking

Although the SCD approach we introduced in Section 4.1 can guarantee the generated API call's faithfulness to the *API documentation*, the faithfulness to the *user's request* is solely dependent on the LLM's zero-shot or few-shot in-context learning and reasoning capabilities. Also, the valuable labeled training data hasn't been exploited yet with SCD. Fine-tuning the LLM might be a straightforward solution, but the huge compute costs associated with the growing parameter size of LLM motivates us to seek alternative solutions.

Inspired by the studies of discriminator-guided generation, we propose the supervised training of a lightweight scorer for the reranking of beamsearched candidate generations. To train the scorer, we generate data as follows: we prompt the Alpaca-13B model with training set samples and obtain associated beam-searched candidate generations for each sample. For each candidate generation, the matching score with respect to the ground truth is calculated as the target of the scorer. Such data is used for the tuning of a RoBERTa-base(Liu et al., 2019) model that has 125M parameters to predict the matching score based on the input of conversation and API Documentation context and the candidate generation. Specifically, we train the model with sample-wise batching that groups candidate generations of the same context into one mini-batch and use the MSE Loss and Spearman Soft Ranking Correlation Loss (Blondel et al., 2020) as the training objective for optimal performance.

374

375

376

377

378

379

380

381

382

385

386

387

388

389

390

391

392

394

395

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

Instead of returning the beam-searched candidate generation that has the highest sequence probability as the final generated API call, we use the trained scorer to discriminate each of the candidate generations and rerank them accordingly. This reranking strategy allows low-cost incorporation of task and domain-specific context reasoning capability learned from the valuable labeled data and compliments the SCD approach and LLM's zeroshot or few-shot in-context learning and reasoning capabilities.

## 5 Experiment Setup

#### 5.1 Datasets

We use the data from **DSTC8** (Kim et al., 2019) and **API Bank** (Li et al., 2023) to conduct the experiments and evaluation of our proposed FAN-TASE framework. Both datasets support the task of API call generation that requires understanding and reasoning of multi-turn human-assistant dialogues. Short-listed APIs and associated Documentation are accompanied by each sample. One major difference is that each DSTC8 sample has two related exemplars while API Bank does not ship with the few-shot in-context learning setting with DSTC8 and the zero-shot in-context learning setting with API Bank. In Appendix A, we provide detailed statistics of DSTC8 and API Bank.

#### 5.2 Baseline and Backbone Models

For in-context learning settings, we include strong baselines GPT3.5-turbo and GPT4 developed by OpenAI, which represents the most recent breakthrough in LLMs with a track record of leading zero-shot and few-shot learning and reasoning capabilities. FANTASE is a plug-and-play modelagnostic approach that can be used in conjunction with any LLM that has an autoregressive decoder producing next-token probabilities. However, the API access of OpenAI models only allows greedy search and does not provide the logits. In consideration of the license, resource constraints, efficiency, and zero-shot/few-shot learning and reasoning capabilities, we opt to use Alpaca-13B as the backbone of our proposed methods. As for baselines of supervised learning settings, we finetune the Alpaca-7B model with DSTC8 training set samples (denoted as AlpDSTC-7B), and we di-

Dataset	DSTC8 (Two Examplars)	API Bank (No Examplars)		
In-context Learning Methods				
Baselines				
GPT4 (est. 1.76T)	37.22	63.66*		
GPT3.5-turbo (est. 175B)	49.28	59.40*		
Alpaca-13B Greedy Search	37.63	24.06		
Alpaca-13B Beam Search	40.49	24.31		
FANTASE with Alpaca-13B as the Base Model				
SCD Greedy Search	42.33	56.64		
SCD Beam Search	44.17	62.66		
Supervised Learning Methods				
	Baselines			
AlpDSTC-7B / Lynx-7B Greedy Search	46.63	48.62		
AlpDSTC-7B / Lynx-7B Beam Search	47.44	50.53		
FANTASE with Alpaca-13B as the Base Model (In-context) and RoBERTa-Base Reranker (Supervised)				
Reranking	46.42	33.33		
SCD Beam Search + Reranking	48.88	64.41		
FANTASE with AlpDSTC-7B / Lynx-7B as the Base Model				
SCD Greedy Search	59.30	65.66		
SCD Beam Search	62.78	67.17		

Table 3: API Call Generation Accuracy Evaluation On DSTC8 and API Bank. For settings involving beam search, we set the beam size to 4. For reproducible results, we set temperature to 0 for all settings. (\* denotes results reported by Li et al. 2023. Best performed In-context and Supervised Learning Methods are **bolded**).

rectly use the Lynx-7B model, an API Bank data tuned Alpaca-7B model, released by Li et al. 2023. We supply detailed information of the aforementioned models in Appendix B.

#### 5.3 Evaluation Settings

423

494

425

426

427

428

429

430

431

432

433

434

435

436

437 438 To verify the effectiveness of FANTASE, we run experiments with the following three evaluation settings that focus on different perspectives:

**API Call Generation Accuracy** measures if the generated API calls fully match their associated ground truth. It is the main metric that reflects if the generation faithfully followed the user's request and the requirements specified in the API documentation. As the order of arguments does *not* matter for both datasets, we calculate unit-wise order-insensitive set matches.

Inference Efficiency measures the time cost of the 439 API call generation. Previous works on constrained 440 441 decoding often vaguely report that the decoding speed is slower than regular decoding algorithms 442 without quantitative measurements. To quantify 443 the speed up brought by FANTASE's State-tracked 444 Constrained Decoding that utilizes CTST (§4.1), 445 we compare the time costs of the API call genera-446 tion with regular/constrained greedy/beam search 447 algorithms using the Alpaca-13B model under the 448 in-context learning setting. 449

450 Context Efficiency measures the effectiveness of
451 our approach in incorporating the API documenta452 tion without the reliance on the repeated supply of
453 the lengthy API documentation in the prompt for
454 in-context learning.

### 6 Results and Analysis

## 6.1 API Call Generation Accuracy

In Table 3, we report the results of API call generation accuracy. The SCD component of FANTASE consistently brings substantial improvements over the base models for both in-context learning settings and supervised learning settings on DSTC8 and API Bank, which demonstrates complementary benefits. 455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

Specifically, for few-shot in-context learning settings evaluated with DSTC8, SCD Greedy Search and SCD Beam Search improve the accuracy by +4.7 and +3.68 over the respective counterparts. For zero-shot in-context learning settings evaluated with API Bank, SCD Greedy Search and SCD Beam Search boost the accuracy by +32.33 and +34.34 respectively, which makes the 13B model's zero-shot generation performance comparable to the GPT3.5-turbo model that has an estimated parameter size of 175B and surpasses the accuracy of fine-tuned 7B model by a large margin. The larger performance gap signifies the value of SCD when labeled data is not available at all.

For supervised learning settings, SCD can still greatly improve the performance of corresponding settings of fine-tuned models and yields +17.04/+16.64 performance gain on API Bank and +12.67/+15.34 performance gain on DSTC8 with greedy/beam search, which leads to the 7B models outperforming GPT models that are 25x~250x times larger in terms of parameter size. Also, we find that GPT4 model's performance on DSTC8 is

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

much worse than the GPT3.5-turbo model. Our diagnosis on the generation reveals that GPT4 model 488 tends to hallucinate the year as 2023 for date-489 related parameters even if the context year is 2019. 490 GPT4 model also made more formatting errors than GPT3.5-turbo on DSTC8. 492

487

491

493

494

495

496

497

498

499

500

504

505

506

507

508

The Reranking component of FANTASE also brings considerable improvements over the base model by supervised training of a lightweight discriminator. The Reranking component alone improves the regular beam search results by +6.15 and +9.02 respectively on DSTC8 and API Bank. The FANTASE framework, with both the SCD component and Reranking component activated, achieves the best overall accuracy showing complementary benefits of the two components. Specifically, for DSTC8, the accuracy of 48.88 is close to the performance of GPT3.5-turbo and better than the supervised fine-tuned model AlpDSTC-7B. For API Bank, the accuracy of 64.41 is better than GPT4 and supervised fine-tuned model Lynx-7B.

#### Inference Efficiency 6.2

	DSTC	DSTC8		API Bank	
Decoding Strategy	Inference Speed (sec/sample)	Speed Up	Inference Speed (sec/sample)	Speed Up	
GS	5.32	-	5.85	-	
SCD GS	3.42	x1.56	3.33	x1.76	
BS	15.12	-	23.15	-	
SCD BS	6.33	x2.39	10.27	x2.25	

Table 4: Generation Speed of Regular Greedy Search (GS) and Beam Search (BS) Decoding versus FAN-TASE's State-tracked Constrained Decoding (SCD) Counterparts.

In Table 4, we quantitatively measure the infer-509 ence time savings of the State-Tracked Constrained 510 Decoding that leverages the CTST as we have in-511 troduced in Section 4.1 and illustrated in Figure 4. 512 Compared to regular greedy and beam search, SCD 513 has the capability of speeding up the API call gen-514 eration by approximately 1.5x~2.4x times. When 515 looking together with generation accuracy, it is 516 quite encouraging that SCD greedy search can 517 achieve regular beam search level's performance 518 with significantly less amount of time and SCD 519 Beam Search can achieve much better performance 520 at regular greedy search level's time cost. For API 521 Bank, SCD Greedy search can even outperform 522 regular bream search with significantly less time cost.

#### 6.3 **Context Efficiency**

Dataset	Setting	w. API Doc	w.o. API Doc	Δ
	GS	37.63	33.74	-3.89
DSTC8	SCD GS	42.33	40.70	-1.63
	BS	40.49	38.24	-2.25
	SCD BS	44.17	42.54	-1.63
	GS	24.06	4.76	-19.3
API Bank	SCD GS	56.64	22.81	-33.83
	BS	24.31	4.51	-19.8
	SCD BS	58.65	23.05	-35.6

Table 5: Generation Accuracy of Regular Greedy Search (GS) and Beam Search (BS) Decoding with / without API Documentation versus FANTASE's State-tracked Constrained Decoding (SCD) Counterparts.

As discussed in Section 4.1, another benefit of SCD is to save the context tokens required in the prompt for supplying the API documentation to the LLM, as SCD is capable of incorporating that information at the decoding stage. Based on the statistics provided in Appendix A, removing API documentation from the input could save an average of 766.28 tokens for DSTC8 and 265.71 tokens for API Bank. In Table 5, we show the performance of SCD when the API documentation is removed from the model's input. For DSTC8, our constrained decoding method manages to maintain the accuracy above 40 when the API Documentation is absent from the input while the unconstrained counterparts suffer larger performance drops of -3.89/-2.25 that leads to larger performance gaps with our approach. For API Bank, the absence of both exemplars and API documentation makes it super challenging for the model to generate relevant API calls as shown by the large performance drop. However, our constrained decoding method can still achieve 22.81/23.05 accuracy in such a scenario, which is close to the performance of the unconstrained version that has API Documentation access in the prompt.

#### 7 Conclusions

The FANTASE framework, with its State-Tracked Constrained Decoding (SCD) and Reranking components, effectively tackles the challenges of generating API calls from complex contexts. By integrating API constraints through a Token Search Trie and employing a lightweight model for reranking, FANTASE not only ensures accurate API call generation but also improves inference and context efficiencies. Its superior performance on the DSTC8 and API Bank datasets confirms FANTASE's significant advancement in enhancing large language models' tool-using ability.

# Limitations

565Despite the effort we have made to the best of our566current ability, we recognize the following limita-567tions of our work:

Evaluation with larger language models: Due to our resource limitations and the input length of the two datasets, the largest models we can fine-tune and infer are 7B and 13B respectively. The effect 571 of the base model's parameter size hasn't been ex-572 tensively studied in this paper, and the impact of our approach on larger models' API call generation 574 accuracy and inference efficiency lacks empirical 575 evidence. Based on the working mechanism of 576 FANTASE, our educated guess is that larger models may make fewer errors which our approach is targeting, so the improvements would be smaller 579 580 than using relatively small models as the backbone. As for the inference efficiency, the absolute time savings should be larger as a forward pass through a larger model takes longer time. However, the relative speed-up would remain at the current level 584 585 as the amount of nodes that only have one child depends on the data instead of the models. To address this issue, we intend to release our code upon the publication of this paper for the ease of further evaluation conducted by other researchers that have sufficient resource.

591Adaptation to a broader range of tasks: FAN-592TASE is specially designed and evaluated for the593task of our interests - API call generation. Although594the high-level idea and concepts of our approach595should be adpatable to other structured text gener-596ation tasks such as SQL generation, table genera-597tion, etc., the adaptation may require considerable598efforts in identifying constraints, signature tokens,599and appropriate constraint-enforcing steps. For un-600structured text generation tasks, it remains unclear601if the identification of signature tokens and appro-602priate constraint-enforcing steps are feasible.

# Ethics Statement

604Our approach significantly enhances the capability605of current models to generate API calls. However,606it's important to acknowledge that the accuracy of607these generations remains imperfect. As API calls608could enable language models to perform tangible609real-world actions, inaccuracies in API call gener-610ation would lead to serious consequences. These611may include but are not limited to: financial losses612when making wrong purchases and reservations,613potential harm to the human being or the environ-

ment when controlling physical objects in unexpected ways, and the dissemination of false or misleading information when retrieving a wrong set of information. Consequently, we urge the users of our methods and the related API call generation models to be aware of such systems' high likelihood of generating inaccurate API calls and to remain vigilant about the possible risks associated with such errors.

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

666

667

668

## References

- Mohaddeseh Bastan, Mihai Surdeanu, and Niranjan Balasubramanian. 2023. NEUROSTRUCTURAL DECODING: Neural text generation with structural constraints. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 9496–9510, Toronto, Canada. Association for Computational Linguistics.
- Mathieu Blondel, Olivier Teboul, Quentin Berthet, and Josip Djolonga. 2020. Fast differentiable sorting and ranking. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
- Xiang Chen, Zhixian Yang, and Xiaojun Wan. 2022. Relation-constrained decoding for text generation. In *Advances in Neural Information Processing Systems*, volume 35, pages 26804–26819. Curran Associates, Inc.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. Plug and play language models: A simple approach to controlled text generation. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages

725

726

728

- 669 670
- 671
- 674 675
- 678
- 679 681
- 682
- 686

- 697
- 701 702 703
- 704
- 706 707 710
- 711
- 712 713
- 714
- 715
- 716 717

718 719

720 721

4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

- 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.
- Chris Hokamp and Qun Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1535–1546, Vancouver, Canada. Association for Computational Linguistics.
- Seokhwan Kim, Michel Galley, Chulaka Gunasekara, Sungjin Lee, Adam Atkinson, Baolin Peng, Hannes Schulz, Jianfeng Gao, Jinchao Li, Mahmoud Adada, Minlie Huang, Luis Lastras, Jonathan K. Kummerfeld, Walter S. Lasecki, Chiori Hori, Anoop Cherian, Tim K. Marks, Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, and Raghav Gupta. 2019. The eighth dialog system technology challenge.
- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. 2021. GeDi: Generative discriminator guided sequence generation. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 4929-4952, Punta Cana, Dominican Republic. Association for Computational Linguistics.
  - Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. 2022. Finetuning can distort pretrained features and underperform out-of-distribution. In International Conference on Learning Representations.
- David Landsman, Jerry Zikun Chen, and Hussain Zaidi. 2022. BeamR: Beam reweighing with attribute discriminators for controllable text generation. In Findings of the Association for Computational Linguistics: AACL-IJCNLP 2022, pages 422-437, Online only. Association for Computational Linguistics.
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023. Api-bank: A comprehensive benchmark for tool-augmented llms.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. CoRR, abs/1907.11692.
- Ximing Lu, Sean Welleck, Peter West, Liwei Jiang, Jungo Kasai, Daniel Khashabi, Ronan Le Bras, Lianhui Qin, Youngjae Yu, Rowan Zellers, Noah A. Smith, and Yejin Choi. 2022. NeuroLogic a\*esque decoding:

Constrained text generation with lookahead heuristics. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 780-799, Seattle, United States. Association for Computational Linguistics.

- Ximing Lu, Peter West, Rowan Zellers, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Neuro-Logic decoding: (un)supervised neural text generation with predicate logic constraints. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4288-4299, Online. Association for Computational Linguistics.
- Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, et al. 2023. Augmented language models: a survey. arXiv preprint arXiv:2302.07842.
- Ansong Ni, Srini Iyer, Dragomir Radev, Ves Stoyanov, Wen-tau Yih, Sida I Wang, and Xi Victoria Lin. 2023. Lever: Learning to verify language-to-code generation with execution. In Proceedings of the 40th International Conference on Machine Learning (ICML'23).
- OpenAI. 2023. GPT-4 technical report. CoRR. abs/2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In NeurIPS.
- Matt Post and David Vilar. 2018. Fast lexically constrained decoding with dynamic beam allocation for neural machine translation. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1314-1324, New Orleans, Louisiana. Association for Computational Linguistics.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. ArXiv, abs/2302.04761.

Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 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, Dominican Republic. Association for Computational Linguistics.

781

783

788

790

791

793

794

796

797

799

809

810

811

812 813

814

815

816

817

818

819

820

821

823

824

825

826

827

830

831 832

833

835

836

- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford\_alpaca.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models.
  - Shufan Wang, Sébastien Jean, Sailik Sengupta, James Gung, Nikolaos Pappas, and Yi Zhang. 2023a. Measuring and mitigating constraint violations of incontext learning for utterance-to-API semantic parsing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7196–7207, Singapore. Association for Computational Linguistics.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. Self-instruct: Aligning language models with self-generated instructions. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.
- Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond.
- Kevin Yang and Dan Klein. 2021. FUDGE: Controlled text generation with future discriminators. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3511–3535, Online. Association for Computational Linguistics.

### Appendix

**A** Datasets

Deteset	DSTC8	API Bank	
Dataset	median / mean	median / mean	
Total Input	1,683	492	
Length (Tokens)	1,644.28	542.86	
API Documentation	735	244	
Length (Tokens)	766.28	265.71	
Examplars	580	N/A	
Length (Tokens)	558.00		
Conversation	200.5	138	
Length (Tokens)	221.08	175.15	
Conversation	6	3	
Turns	6.31	3.21	
Target API Call	46	27	
Length (Tokens)	45.70	35.17	
Target API Call	3	2	
Arguments Amount	3.42	2.27	

Table 6: Statistics of DSTC8 and API Bank Data.

In Table 6, we provide the detailed statistics of DSTC8 and API Bank. The inputs and outputs of DSTC8 are longer than API Bank counterparts. DSTC8 also has more turns of conversation and API call arguments than API Bank, which makes DSTC8 a more challenging dataset in terms of reasoning the context and generating appropriate API calls even if two exemplars are provided for each sample.

Our experiments and evaluation of the FAN-TASE framework are based on the test split of DSTC8 and API Bank that has 490 and 399 samples respectively. For DSTC8, we remove the longest one out of the 490 samples as it causes CUDA out of memory error on our server. 857

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

871

872

858

## **B** Models

873

Alpaca-13B is a LLaMA-based instruction following model that has 40 layers, a hidden size of 5120,
40 self-attention heads and 13 billion parameters.
For our experiments, we use Huggingface checkpoint chavinlo/gpt4-x-alpaca.

AlpDSTC-7B and Lynx-7B are the Alpaca-7Bbased model that has 32 layers, a hidden size of 4096, 32 self-attention heads and 7 billion parameters. Lynx-7B is the training set fine-tuned model released by API Bank authors (Huggingface checkpoint: liminghao1630/Lynx-7b). The model was tuned for 3 epochs with a learning rate of 2e-5 and an effective batch size of 256. We follow the same setting of Lynx-7B to tune a Alpaca-7B with DSTC8 training set samples.

RoBERTa-base is a BERT-based model that has 12
layers, a hidden size of 768, 12 self-attention heads,
and 125 million parameters. For our experiments,
we use Huggingface checkpoint roberta-base
and tune the model with training set candiate generations for 5 epochs with a learning rate of 5e-5
and an effective batch size of 256.

GPT3.5 and GPT4 are only accessible via API
or web interface. Details of these two models
have not been officially released by OpenAI, but a
broadly accepted latency-based parameter size estimation is 175 billion parameters for GPT3.5 and
1.76 trillion parameters for GPT4. For our experiments, we use checkpoints gpt-3.5-turbo-0613
and gpt-4-0613.

904 License information of Alpaca can be
905 found at https://github.com/tatsu-lab/
906 stanford\_alpaca/blob/main/LICENSE and
907 https://github.com/tatsu-lab/stanford\_
908 alpaca/blob/main/DATA\_LICENSE.

- 909License information of LLaMA can be found910at https://docs.google.com/forms/d/e/9111FAIpQLSfqNECQnMkycAp2jP4Z9TFX0cGR4uf7b\_912fBxjY\_OjhJIL1KGA/viewform.
- 913 License information of Lynx-7B and API Bank
  914 data can be found at https://github.com/
  915 AlibabaResearch/DAMO-ConvAI/blob/main/
  916 api-bank/LICENSE.
- License information of DSTC8 data 917 be found at 918 can https://github. com/google-research-datasets/ 919 920 dstc8-schema-guided-dialogue/blob/
- 921 master/LICENSE.txt.
- 922License information of RoBERTa can be found923at https://github.com/facebookresearch/

# fairseq/blob/main/LICENSE.

License information of GPT-X models can 925 be found at https://openai.com/policies/ 926 terms-of-use. 927