

The Whole Truth and Nothing But the Truth: Faithful and Controllable Dialogue Response Generation with Dataflow Transduction and Constrained Decoding

Anonymous ACL submission

Abstract

In a real-world dialogue system, generated text must satisfy several interlocking constraints: informativeness, truthfulness, and ease of control. The two predominant paradigms in language generation—neural language modeling and rule-based generation—struggle to satisfy these constraints simultaneously. We describe a hybrid architecture for dialogue response generation that combines the strengths of both paradigms. The first component of this architecture is a rule-based content selection model defined using a new formal framework called *dataflow transduction*, which uses declarative rules to transduce a dialogue agent’s actions and their results (represented as dataflow graphs) into context-free grammars representing the space of contextually acceptable responses. The second component is a constrained decoding procedure that uses these grammars to constrain the output of a neural language model, which selects fluent utterances. Our experiments show that this system outperforms both rule-based and learned approaches in human evaluations of fluency, relevance, and truthfulness.

1 Introduction

In a task-oriented dialogue system, response generation is naturally posed as a conditional language modeling problem: dialogue agents must produce a contextually appropriate natural language string conditioned on the history of the user and agent interaction. But unlike many language generation problems, a good dialogue response generation model is not (just) a model of typical human utterances in context. Instead, effective dialogue agents must balance fluent generation with a set of much stricter constraints.

Consider the dialogue shown in Fig. 1. In the first turn of this dialogue, the user makes a request, the dialogue agent correctly translates it into a computation—here represented as a dataflow

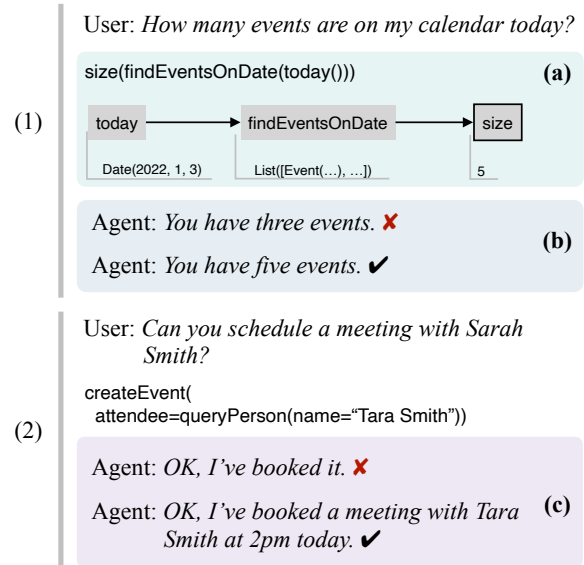


Figure 1: Interaction between a user and a dialogue agent. Once the user’s request is translated into an agent action—expressible as a program or dataflow graph (a)—the agent must generate a response. Agent responses might simply state the result of the agent’s action, but must do so truthfully (b). Often responses should describe both the action and the result, *e.g.*, to help users identify when the agent has misunderstood their request (c). These responses should be straightforward for system designers to inspect and modify.

graph (Fig. 1a)—then it needs to accurately describe this computation’s return value (Fig. 1b), rather than using an arbitrary number (*e.g.*, 3 in the Date) on the dataflow graph. In the second step, the agent may also make a mistake: perhaps because of a speech recognition error, it creates a meeting with *Tara Smith* rather than *Sarah Smith*. Simply describing the result of its action might cause a user to incorrectly conclude that their request was completed successfully. To avoid confusion, a system designer might wish to ensure that the agent instead echoes back to the user the details of the agent’s action (Fig. 1c). This example highlights the challenges central to building real-world dialogue response generation systems.

057 First, response generation is not simply a prob- 108
058 lem of describing the *result* of a computation in 109
059 natural language. In some cases, response gener- 110
060 ators may also usefully **describe the provenance** 111
061 of that result—the computation itself and its in- 112
062 termediate values. In many human-to-human con- 113
063 versations, a response as detailed as Fig. 1c would 114
064 be over-informative, violating Grice’s maxim of 115
065 quantity (1975). But for a speaker that is prone 116
066 to mistakes, such as an AI agent, describing its 117
067 own understanding can increase user trust when 118
068 the understanding is accurate and provides an op-
069 portunity for correction when it is not.

070 Second, dialogue response generation systems 119
071 must **guarantee truthfulness**: as typically the pri- 120
072 mary source of information about the action that a 121
073 dialogue agent took, a response generator that de- 122
074 scribes even a small fraction of these computations 123
075 incorrectly can produce disastrous results. Import- 124
076 antly, truthful utterances might be low-probability 125
077 under a domain-general language model (LM), 126
078 particularly when they reflect errors in language 127
079 understanding (as in Fig. 1b). 128

080 Finally, response generation systems must **sup- 129
081 port declarative specification of agent behavior.** 130
082 When confusing or infelicitous responses are dis- 131
083 covered, it should be possible to easily and pre- 132
084 cisely modify them without changing the dialogue 133
085 agent’s behavior in other contexts. 134

086 In recent years, the main focus of academic di- 135
087 alogue research has been on “end-to-end” learned 136
088 models for response generation, especially neural 137
089 sequence models (Vinyals and Le, 2015; Zhang 138
090 et al., 2020b). But while such models excel at pro- 139
091 ducing fluent and coherent output, research con- 140
092 tinues to find that they struggle in maintaining 141
093 faithfulness (Wiseman et al., 2017; Maynez et al., 142
094 2020). Perhaps more fundamentally, because the 143
095 behavior of such systems is encoded implicitly in 144
096 their training data, designing a dialogue system 145
097 requires system builders to write and edit a large 146
098 number of training examples whose final effect 147
099 may be difficult to predict.

100 As a result, many dialogue systems in the real 148
101 world remain rule-based: system builders hand- 149
102 write rules (*e.g.*, in the form of a synchronous 150
103 grammar) for transforming dialogue states into 151
104 text, and these rules are applied directly during de- 152
105 ployment. But such rule-based systems are also 153
106 notoriously difficult to build and maintain (Walker 154
107 et al., 2002; Reiter, 2022). They require designers

108 to anticipate every low-level question about sur- 109
110 face realization, and to encode these in the same 110
111 grammar that is responsible for enforcing high- 111
112 level properties like truthfulness. 112

113 Given the many strengths of modern LMs, is 113
114 there a way to leverage them while satisfying the 114
115 numerous other demands on dialogue response 115
116 generation systems? In this paper, we describe a 116
117 hybrid approach that combines the advantages of 117
118 end-to-end and rule-based approaches. This ap-
118 proach has two components:

- 119 • A dataflow transduction procedure, based on 119
120 a new formalism that uses declarative rules to 120
121 map a computation (represented as a dataflow 121
122 graph) into a context-free grammar (CFG) 122
123 that defines the space of all responses al- 123
124 lowed for the given computation. This formal 124
125 framework makes it possible to write rules 125
126 to precisely and truthfully describe both data 126
127 and its provenance, while performing supple- 127
128 mentary computation where needed to pro- 128
129 duce informative responses. 129
- 130 • A constrained decoding procedure that inter- 130
131 sects a CFG with a neural LM, making it pos- 131
132 sible to decompose language generation into 132
133 a **content selection model** (implemented by 133
134 the grammar) and a separate **fluency model** 134
135 (implemented by an LM). 135

136 Together, dataflow transduction and constrained 136
137 decoding make it possible to build a faithful gen- 137
138 eration system capable of describing a complex, 138
139 open-ended, space of tasks. Using a subset of 139
140 SMCaFlow dialogues (Semantic Machines et al., 140
141 2020) and only 187 declarative rules, our hybrid 141
142 system is consistently rated as more truthful, rele- 142
143 vant, and fluent than either a rule-based or end-to- 143
144 end neural system. Similar results are observed on 144
145 MultiWOZ dialogues (Budzianowski et al., 2018; 145
146 Eric et al., 2020). Code, data, and trained models 146
147 used in our experiments will be released. 147

148 2 Problem Formulation

149 We study the problem of response generation for 149
150 task-oriented dialogue. A dialogue, like the one in 150
151 Fig. 1, consists of a sequence of **turns**, each 151
152 consisting of a **user utterance** x_i , one or more **actions** 152
153 a_i , and an **agent response** y_i . The job of a **dia-** 153
154 **logue agent** is to predict an appropriate action and 154
155 response from a dialogue history, *i.e.*, mapping 155
156 from $(x_1, a_1, y_1, x_2, a_2, y_2, \dots, x_n) \mapsto (a_n, y_n)$. 156

A common approach to building dialogue agents decomposes this prediction process into several steps. First, a **language understanding module** maps from a user utterance (and possibly other components of the dialogue history) to a meaning representation (e.g., a structured user intent, API request or executable program). This meaning representation is evaluated, producing actions a , which are passed to a **response generation module** that produces an agent utterance y .

The focus of this paper is the response generator. We assume that we have a pre-specified language understanding module that maps from conversation histories to **computations**, in the form of short programs, which are then executed to produce actions a . As described by [Semantic Machines et al. \(2020\)](#), these computations may equivalently be viewed as **dataflow graphs** in which each node is labeled with a function, constructor, or primitive value, as well as a return value once the node is executed. We additionally assume access to a dataset of dialogues containing gold-standard user and agent utterances. Given a language understanding module and a dataset of dialogues, we aim to implement a response generator that, when applied to a dataflow graph, satisfies the three properties outlined in §1: description of data and its provenance, guaranteed truthfulness, and declarative specification.

Our response generation system is built from two pieces: (1) a procedure for transducing dataflow graphs into CFGs (§3), and (2) a constrained decoding procedure for intersecting a CFG with a neural LM (§4). Hybrid generation systems of this kind have a long history in NLP ([Langkilde and Knight, 1998](#)). Our aim in this paper is to show the benefits of a new generation paradigm based on dataflow transduction, and offer new rule-writing formalisms and decoding algorithms tailored to modern language models.

3 Dataflow Transduction

Given a dataflow graph G (e.g., Fig. 1a) rooted at a node v_{root} (the return value of the program represented by the dataflow graph), our task is to generate a string that describes v_{root} and its provenance. To achieve this, we propose a new formal framework for generation based on **dataflow transduction**. At a high level, the formalism uses declarative rules that describe how to transform a dataflow graph into a graph-specific gram-

Head: S
Body:
<pre>match computation: case findEventsOnDate(date): num = size(computation) event = head(computation) return {"num": num, "event": event, "date": date}</pre>
Response Template:
<i>I found {LEX <num>} event {PP <date>}. It's {EVENT <event>}.</i>

Figure 2: A dataflow transduction rule with head S, a body (expressed in Python), and a response template (which queries the dictionary returned by the body).

mar (specifically a **quasi-synchronous context-free grammar**, or QCFG) that defines the space of allowed responses. These rules walk along the graph, introduce new computations (dataflow sub-graphs) as needed, and add rules to the grammar.

Formally, a dataflow transducer \mathcal{S} is defined by a 4-tuple $(\mathcal{T}, \Sigma, \mathcal{R}, t_{\text{start}})$ where \mathcal{T} is a set of nonterminal types,¹ Σ is the set of terminals (word types), \mathcal{R} is a set of dataflow transduction rules (see §3.1), and t_{start} is the nonterminal type of the start symbol. When applied to G the dataflow transducer expands the graph, yielding a new graph \bar{G} , and produces a QCFG.

A QCFG ([Smith and Eisner, 2006](#)) is a specialized CFG whose nonterminals include alignments to the nodes $V(\bar{G})$ of \bar{G} . Where an ordinary CFG might specify ways to generate an NP (noun phrase) or a DATE, a QCFG would specify ways to generate an NP or DATE that describes the result and provenance of v , for each appropriately typed node $v \in V(\bar{G})$. A QCFG resulting from dataflow transduction is a 4-tuple $(\mathcal{T} \times V(\bar{G}), \Sigma, \mathcal{P}, t_{\text{start}})$ where $\mathcal{T} \times V(\bar{G})$ is the QCFG's set of nonterminals and \mathcal{P} is its set of productions. A QCFG production has the form $\alpha \rightarrow \beta_1\beta_2 \cdots \beta_N$ where the left-hand-side $\alpha = (t, v) \in \mathcal{T} \times V(\bar{G})$ is a QCFG nonterminal, and each β_i can be either a nonterminal (t_i, v_i) or a terminal in Σ . The v_i of a right-hand-side nonterminal β_i may have appeared in the original G , or may have been added to \bar{G} by the dataflow transducer. These production rules then derive a set of strings as in an ordinary CFG.

3.1 Dataflow Transduction Rules

A dataflow transduction rule is applied to a node $v \in \bar{G}$ (if v has appropriate properties) to create a

¹In practice, nonterminal types might correspond to dialogue acts, syntactic categories, semantic categories, etc. This is up to the designer.

single QCFG production $(t, v) \rightarrow \dots$ that could be used to describe v . An example rule is shown in Fig. 2. A rule has three components: (1) a **head**, namely the nonterminal type $t \in \mathcal{T}$; (2) a **body**, which is a piece of code that determines whether the rule can apply to v , and which may look up or create nodes that are related to v ; and (3) a **response template**, which specifies the right-hand side of the QCFG production in terms of the related nodes that identified in the body.

Rule Head. This nonterminal type characterizes the type of node that the transduction rule is able to describe and the type of description that it will produce.¹ When a rule with head t is successfully applied to the node v , the resulting QCFG production has left-hand-side (t, v) .

Rule Body. A rule body declares the condition when the rule can be applied by examining the dataflow graph \bar{G}_v rooted at v . It can contain executable logic that identifies additional computation nodes that will be recursively described.² For example, the rule body in Fig. 2 checks whether \bar{G}_v has the form `findEventsOnDate(date)`. If so, it binds the variable `date` accordingly, and introduces new nodes into \bar{G} , bound to the variables `num` and `event`, which compute the number of events and the first event. All three of these variables will be referenced in the response template.

Response Template. The response template says how to create the right-hand side of the QCFG rule—a sequence $\beta_1 \dots \beta_N$ of terminals and nonterminals. Each QCFG nonterminal $\beta_i = (t_i, v_i)$ specifies a related node $v_i \in V(\bar{G})$ to describe, along with a dataflow nonterminal t_i that says *how* to describe it. The possible descriptions of v_i will thus emerge from applying transducer rules with head t_i to node v_i . In our template syntax, the notation `{EVENT <event>}` would construct the QCFG nonterminal (EVENT, v) , if the rule body has bound the variable `event` to the node v . This syntax is illustrated in Fig. 2; *e.g.*, the response template will construct three QCFG nonterminals, with types LEX, PP, and EVENT.

3.2 Dataflow Transduction Procedure

Given a dataflow transducer \mathcal{S} and a dataflow graph G rooted at node v_{root} , we can transduce the graph into a QCFG as follows. The system starts out by creating QCFG productions that can

²Note that the nodes added by the body may represent further computations on existing nodes of \bar{G}_v or may be completely disjoint from the existing nodes.

expand the start nonterminal $(t_{\text{start}}, v_{\text{root}})$. For each transduction rule in \mathcal{R} whose head is t_{start} , it executes the body, which checks any additional conditions for whether the rule can be applied to v_{root} , binds variables, and uses the response template to create a QCFG production. If these productions mention new nonterminals, the system recursively creates further QCFG productions, in the same way, that can expand those nonterminals. As a special case, to expand a nonterminal of the form (LEX, v) , the system creates a QCFG production whose right-hand side gives the value of v , as rendered into natural language using a lexicalization function rather than a template; *e.g.*, a value `Long(1)` would be rendered as “1”.

The recursive process continues until productions have been created for every nonterminal that appears in the QCFG. The resulting QCFG compactly represents a combinatorial space of possible responses. It will generally include multiple productions aligned to the same node v , created by different dataflow transduction rules.

This mechanism can be used to copy simple values like strings and numbers from the dataflow graph, as well as to create more complex recursive descriptions. Note that (1) transduction rules are selected via their head but also condition on the dataflow graph through their body, and (2) all QCFG nonterminals are grounded in the dataflow graph. Together, this provides a means to ensure truthfulness when generating responses.

4 Constrained Decoding

In this section, we describe how to integrate the formal framework above with a general LM to perform response generation, as illustrated in Fig. 3. Given a derived QCFG of the kind described in §3.2, we perform constrained decoding as in (Shin et al., 2021; Roy et al., 2022), generating response candidates from a pretrained LM.

The QCFG resulting from dataflow transduction implicitly represents a set of possible derivation trees and the agent responses they yield. As long as transduction rules faithfully describe the nodes they apply to, every derivation in this set will correspond to a truthful agent utterance. But these utterances may not always be grammatical or natural. For example, the response template in Fig. 2 may be realized as “I found 2 event on Monday” since the rule body does not check whether the value of `num` is 1. Similarly, the response template

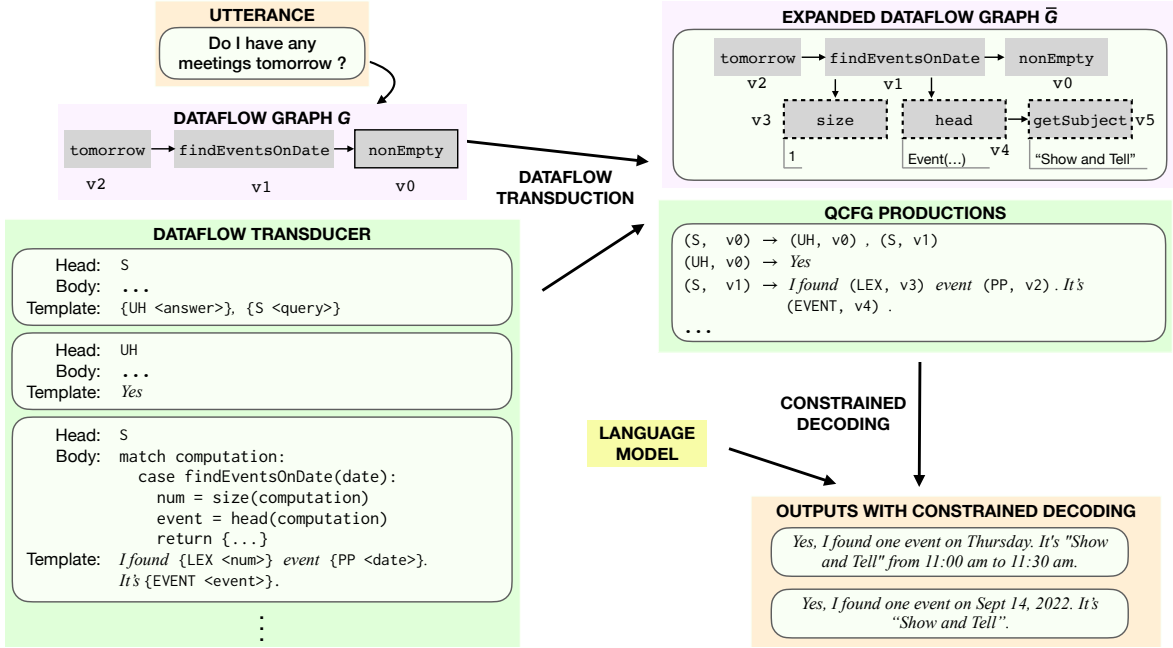


Figure 3: The hybrid response generation approach using dataflow transduction and constrained decoding. Given a computation `nonEmpty(findEventsOnDate(tomorrow()))` for the user utterance “Do I have any meetings tomorrow”, we first derive QCFG productions by applying the dataflow transducer to the dataflow graph G using the procedure described in §3.2. This procedure also expands the dataflow graph into \bar{G} : for example, the nodes `v3` and `v4` were added by the third transducer rule. Then we extract candidate responses from a LM, constrained by the QCFG. The varying descriptions of the date `v2` and the event `v4` are permitted because the QCFG offers a choice of productions that can be used to expand the $(PP, v2)$ and $(EVENT, v4)$ nonterminals. (Those productions and the transducer rules that created them are not shown in the figure. The nodes added by those transducer rules and used by those productions are also not shown, except for `v5`.)

{EVENT $\langle event \rangle$ } starts on {DATE $\langle date \rangle$ }.

may be realized as *The product meeting on Monday starts on Monday*, if the grammar permits identifying events by their dates. With carefully engineered and highly specialized rules (e.g., using extremely fine-grained nonterminal types), it would be possible to ensure that the responses are always fluent and even that there is always a single possible outcome from the top-down search procedure. However, this would usually require much a more complicated set of rules, which creates a burden for system development and maintenance.

Our proposed approach instead uses a large-scale pretrained LM (preferably fine-tuned) to select among truthful utterances produced by the QCFG.³ One option is to use the LM to re-rank all strings that can be produced by the QCFG, but that would be very computationally expensive. Instead, we follow Shin et al. (2021) and Roy et al.

³Of course, decisions deferred to the LM could be encoded in the grammar instead. While this is rarely necessary to ensure grammaticality or fluency, system designers might choose to encode some *pragmatic* decisions, like how much detail to provide, in the grammar rather than in the LM.

(2022), who decode sentences from a given LM under the constraint that they must be valid under a given CFG. This constrained decoding method uses Earley’s algorithm (1970) to incrementally parse the sentence as it is generated and determine the set of words that could grammatically serve as the next token. In contrast to these prior papers, which used a static CFG, we derive a new CFG each time the dialogue agent needs to generate a response, by applying the dataflow transducer to the current dataflow graph.

5 Experiments

To evaluate this approach, we conducted a set of detailed experiments on the SMCaFlow dataset (Semantic Machines et al., 2020) (§5.1–§5.3), and a brief study on applying our approach to the MultiWOZ dataset (Budzianowski et al., 2018) (§5.4).

5.1 Data and Evaluation Metrics

SMCaFlow is a large-scale task-oriented dialogue dataset, in which each user utterance is annotated with a correct dataflow program (i.e., computa-

System	Automatic Metrics					Human Evaluation (%)		
	BLEU	ROUGE	BERTSc.	R@1	R@5	Grammatical	Relevant	Truthful
QCFG Random Sampling	.35	.58	.50	.02	.06	62.3 [†]	90.9 [†]	92.3
Unconstrained Decoding	.77	.86	.87	.47	.66	98.7	93.3	82.2 [†]
QCFG-Constrained Decoding	.80	.87	.86	.56	.78	99.0	96.6	91.6
Gold	1.0	1.0	1.0	1.0	1.0	99.0	98.0	92.3

Table 1: Evaluation results on SMCaFlow. Automatic metrics are calculated against the gold responses on the full validation set. Human evaluation is conducted on 297 randomly sampled validation examples. †: Results are significantly worse than the “Gold” system ($p < 10^{-4}$, McNemar’s test).

tion) and a “gold” response that would be desirable for the agent to produce.⁴ We use the v2.0 release processed by Platanios et al. (2021). We focus on a subset of SMCaFlow involving calendar event queries. This subset contains 8938 training examples and 1041 validation examples. We found that 187 transduction rules, written by some of us in a matter of hours, were sufficient to cover all gold system responses in these examples.⁵

Automatic Metrics. For automatic evaluation, we use several reference-based metrics, including BLEU-4 (Papineni et al., 2002), ROUGE-L (Lin, 2004), and BERTScore-F1 (Zhang et al., 2020a), computed using the GEM-metrics tool.⁶ Following the recommendation in Zhang et al. (2020a), we use the re-scaled version of BERTScore which is easier to interpret. We additionally consider exact match scores, i.e., **R@K**, which measure whether one of the top K response candidates exactly matches the reference. Both **R@1** and **R@5** scores are reported. We lowercase all the strings and remove any extra spaces while computing the exact match between two strings.

Human Evaluation. It is well-known that popular automatic evaluation metrics may not always reflect the true quality of the generated responses (Celikyilmaz et al., 2021). Thus, we further carry out human evaluation on 297 examples randomly sampled from the validation data. Specifically, for each generated response, we collect human judgments on three questions: **grammaticality** (“has the virtual assistant made any grammar errors?”), **relevance** (“has the virtual assistant mis-

understood the user’s request?”), and **truthfulness** (“has the virtual assistant provided any incorrect information as judged using the database and timestamp?”). Three judgments are collected for each question, and we report the percentage of examples where “no” is the majority-voted answer. Higher percentages are better. Crowdworkers are recruited from Amazon Mechanical Turk with qualification requirements such as having a work approval rate higher than 80% and having performed a minimum of 100 annotations. They are paid at the rate of \$0.15 per judgment. For responses generated by the constrained decoding approach, the inter-annotator agreements for the three questions are around 90%, 78% and 76%, respectively, as measured by the percentage of examples where all three workers choose the same answer. More details are provided in Appendix A.

5.2 Main Results

Our main evaluation results on SMCaFlow are shown in Table 1. The first baseline we considered is to randomly sample responses from the generated QCFG. The other baseline is unconstrained LM decoding without using dataflow transduction. Model outputs are compared to human-authored agent utterances. For both unconstrained and constrained decoding, we prompt the LM with a string representation of the computation graph (i.e., the format released in SMCaFlow v2.0), followed by its execution result rendered as a JSON string. We use beam search with a beam size $K = 5$. The LM is initialized from CodeT5-base (Wang et al., 2021) and fine-tuned on all training examples. See Appendix B for more details.

As expected, the QCFG random sampling baseline struggles on all the automatic metrics, since dataflow transduction rules are written with an emphasis on truthfulness rather than fluency. This is reflected in the grammaticality score from the

⁴The “gold” responses are generated from a production system that includes rule-based constraints and manually validated by human experts, according to the dataset authors.

⁵Some of our rule bodies chose to expand the dataflow graph by calling functions, so we also had to implement those functions. In an end-to-end dialogue system, most of those functions would already have been implemented to support agent actions, not just natural language responses.

⁶<https://github.com/GEM-benchmark/GEM-metrics>

human evaluation as well. However, the truthfulness score is as high as 92.3%, indicating the generated responses are rarely incorrect. Generated responses are sometimes generic or omit relevant information for the user request, which partially contributes to the high truthfulness score, but is reflected in the relevance score, which is the lowest among all compared approaches.

In contrast, unconstrained decoding without dataflow transduction achieves impressive scores on automatic evaluation. Human evaluation also suggests that the generated responses are grammatically correct and relevant to the user’s request in most cases. However, unconstrained decoding scores low on truthfulness, making false statements in about one-fifth of the generated responses. This high rate of factual errors from neural LMs is consistent with findings in prior work (Wiseman et al., 2017; Maynez et al., 2020). It is usually unacceptable in real-world applications.

Compared with unconstrained decoding, our proposed QCFG-constrained decoding achieves significantly better scores on exact match, truthfulness, and even relevance, while maintaining similar scores on BLEU, ROUGE, BERTScore and grammaticality. In particular, human evaluation results indicate that the quality of generated responses is very close to that of the gold responses. We share some qualitative analysis in Appendix C.

Since even the gold responses did not achieve 100% on human evaluation scores, we manually inspected those problematic examples. There are 4 examples for which the majority-voted answer to the ungrammaticality question is “*yes but understandable*”, and others are all rated as not containing any grammar errors. For the relevance question, 4 examples are due to arguably bad data and 2 examples receive tied votes. For the truthfulness question, 9 examples are due to arguably bad data, 8 examples are due to crowd worker mistakes, and 6 examples receive tied votes.

5.3 Ablation Study

We next analyze how the amount of fine-tuning data and the context used in the input sequence impact the quality of generated responses. Results are summarized in Table 2.

Impact of fine-tuning: Without fine-tuning the LM, neither unconstrained nor constrained decoding works well. This is likely due to the mismatch between the pre-training tasks and the re-

	BLEU	ROUGE	BERTSc.	R@1	R@5
1. LM without fine-tuning					
✗	.45	.03	-.29	.00	.00
✓	.38	.22	.07	.02	.02
2. LM fine-tuned on 3% training data					
✗	.68	.82	.80	.26	.40
✓	.73	.83	.81	.39	.61
3. LM fine-tuned on full training data					
✗	.77	.86	.87	.47	.66
✓	.80	.87	.86	.56	.78
4. LM input without execution results					
✗	.58	.70	.72	.27	.42
✓	.78	.86	.84	.54	.77
5. LM input with user utterance					
✗	.76	.88	.87	.45	.65
✓	.77	.85	.85	.54	.78

Table 2: SMCaFlow ablation results, varying the amount of fine-tuning data (groups 1–3) and the context used in the input sequence (groups 4–5). ✗ and ✓ on the first column use unconstrained and QCFG-constrained decoding, respectively.

sponse generation task. However, after fine-tuning on only a random 3% of the training data, both approaches achieve significantly better scores, with larger gains on QCFG-constrained decoding. This suggests that QCFG-constrained decoding is much more data-efficient in the low-data regime. Indeed, using 3% of the training data, QCFG-constrained decoding is on par with the unconstrained decoding with 100% of the training data, indicating that several expert hours spent on creating dataflow transduction rules can dramatically reduce the cost of collecting training data. While gaps between unconstrained and QCFG-constrained decoding on automated metrics are small in the full-data setting (Table 1), unconstrained decoding still performs poorly on the truthfulness evaluation. Thus, truthfulness failures from unconstrained decoding are not straightforwardly solved by scaling up training data; QCFG-constrained decoding offers an easier path to faithful response generation.

Impact of context: Results in groups 3–5 in Table 2 all use the full training data to fine-tune the LM. The difference is in the context used in the input sequence to the LM. For group 3, the input sequence is the computation concatenated with the execution result, which is the same setup used in §5.2. For group 4, we omit the execution results from the LM input (but not from the decoder

constraints), whereas for group 5, we add the user utterance (prefixed to the computation). Comparing group 3 and group 4, omitting execution results significantly harms the performance of unconstrained decoding. In contrast, dataflow transduction rules can execute the computation internally, and do not require the LM to condition on it. Comparing group 3 and group 5, adding user utterances to prompts does not bring any additional benefits to both approaches.

5.4 Experiments with MultiWOZ Dataset

To demonstrate the general applicability of our approach for response generation, we carry out a brief study on the widely used MultiWOZ 2.1 dataset (Budzianowski et al., 2018; Eric et al., 2020). We automatically convert the system act annotations to dataflow computations and write 14 transduction rules. For generating responses, we use the predicted system acts from the MT-TOD system (Lee, 2021). Similar to our experiments on SMCaFlow, we fine-tune CodeT5-base on all training examples, using the ground-truth belief state and predicted system act as the input sequence. For evaluation, we randomly sample 100 examples from the test split, and two authors manually rate the generated responses from our QCFG-constrained decoding system and the MT-TOD system. The agreement is 100%. Almost all generated responses are grammatically correct and relevant to the user utterance. To rate truthfulness, we use the predicted system acts as the references. Our QCFG-constrained decoding approach produce truthful responses for all 100 examples, whereas only 89 responses from the MT-TOD system are truthful with respect to its predicted system act. Among the 11 remaining examples, 7 of them are due to imperfect delexicalization and 4 are due to hallucination.

6 Related Work

One line of response generation research focuses on generating fluent and coherent responses directly from user utterances without any intermediate structured representation. This paradigm is mostly used for chatbots, as in early rule-based systems (Weizenbaum, 1966; Wallace, 2009), neural conversation models (Vinyals and Le, 2015; Shang et al., 2015; Sordani et al., 2015; Li et al., 2016; Serban et al., 2016), and recent large-scale pretrained LMs like DialoGPT (Zhang et al., 2020b) and GPT-3 (Brown et al., 2020).

Another line focuses on generating text from structured data, with applications beyond dialogue response generation. For example, the WebNLG challenge (Gardent et al., 2017) generates natural language descriptions from relation tuples, and Lebrete et al. (2016) generate a biography from a structured “infobox” record. Many recent dialogue response generation tasks adopt dialogue-act-based meaning representations, including the MultiWOZ dataset (Budzianowski et al., 2018), the Schema-Guided dialogue dataset (Rastogi et al., 2020), and the E2E NLG challenge (Dusek et al., 2020). In contrast, our response generation task uses computations as the input, which do not directly encode the dialogue acts of the responses. This is a more challenging task, as the system needs to perform extra reasoning to obtain the derived information. In this sense, our task is similar to the one in CoSQL (Yu et al., 2019) and Logic2Text (Chen et al., 2020).

Constrained decoding techniques for neural LMs have been developed for text generation with different types of constraints (Balakrishnan et al., 2019; Dathathri et al., 2020; Lu et al., 2021, 2022). Shin et al. (2021) develop a constrained decoding approach for semantic parsing by restricting the LM output at each step according to a given grammar. Differently, the grammar productions in our case are derived dynamically for each input.

7 Conclusion

We have described a hybrid approach for building dialogue response generation systems. Our approach introduces a new formalism for transducing a dataflow graph into a QCFG, which is then used in a constrained decoder that intersects the QCFG with a neural LM. (In future work, the QCFG could be weighted to express its own preferences.) This formal framework makes it possible to write rules to precisely and truthfully describe data and its provenance while deferring surface realization decisions to a flexible language model.

This new approach outperforms unconstrained conditional language modeling in both automatic and human evaluations, especially on truthfulness. Moreover, using 3% of the training data, the constrained decoding approach is on par with the unconstrained decoding approach when it uses 100% of the training data, indicating that several expert hours spent on authoring rules can dramatically reduce the cost of data annotation.

References

- 631 Anusha Balakrishnan, Jinfeng Rao, Kartikeya Upasani,
632 Michael White, and Rajen Subba. 2019. **Con-**
633 **strained decoding for neural NLG from composi-**
634 **tional representations in task-oriented dialogue.** In
635 *Proceedings of the 57th Conference of the Associa-*
636 *tion for Computational Linguistics, ACL 2019, Flo-*
637 *rence, Italy, July 28- August 2, 2019, Volume 1:*
638 *Long Papers*, pages 831–844. Association for Com-
639 putational Linguistics.
- 640 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
641 Subbiah, Jared D Kaplan, Prafulla Dhariwal,
642 Arvind Neelakantan, Pranav Shyam, Girish Sastry,
643 Amanda Askell, Sandhini Agarwal, Ariel Herbert-
644 Voss, Gretchen Krueger, Tom Henighan, Rewon
645 Child, Aditya Ramesh, Daniel Ziegler, Jeffrey
646 Wu, Clemens Winter, Chris Hesse, Mark Chen,
647 Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin
648 Chess, Jack Clark, Christopher Berner, Sam Mc-
649 Candlish, Alec Radford, Ilya Sutskever, and Dario
650 Amodei. 2020. **Language models are few-shot**
651 **learners.** In *Proceedings of Advances in Neural*
652 *Information Processing Systems*, volume 33, pages
653 1877–1901. Curran Associates, Inc.
- 654 Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang
655 Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ra-
656 madaan, and Milica Gašić. 2018. **MultiWOZ - a**
657 **large-scale multi-domain Wizard-of-Oz dataset for**
658 **task-oriented dialogue modelling.** In *Proceedings of*
659 *the 2018 Conference on Empirical Methods in Natu-*
660 *ral Language Processing*, pages 5016–5026, Brus-
661 sels, Belgium. Association for Computational Lin-
662 guistics.
- 663 Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao.
664 2021. **Evaluation of text generation: A survey.**
665 *arXiv:2006.14799v2 [cs.CL]*.
- 666 Zhiyu Chen, Wenhui Chen, Hanwen Zha, Xiyou
667 Zhou, Yunkai Zhang, Sairam Sundaresan, and
668 William Yang Wang. 2020. **Logic2Text: High-**
669 **fidelity natural language generation from logical**
670 **forms.** In *Findings of the Association for Computa-*
671 *tional Linguistics: EMNLP 2020*, pages 2096–2111,
672 Online. Association for Computational Linguistics.
- 673 Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane
674 Hung, Eric Frank, Piero Molino, Jason Yosinski, and
675 Rosanne Liu. 2020. **Plug and play language mod-**
676 **els: A simple approach to controlled text genera-**
677 **tion.** In *Proceedings of 8th International Confer-*
678 *ence on Learning Representations, ICLR 2020, Ad-*
679 *dis Ababa, Ethiopia, April 26-30, 2020*. OpenRe-
680 view.net.
- 681 Ondrej Dusek, Jekaterina Novikova, and Verena Rieser.
682 2020. **Evaluating the state-of-the-art of end-to-end**
683 **natural language generation: The E2E NLG chal-**
684 **lenge.** *Computer Speech and Language*, 59:123–
685 156.
- Jay Earley. 1970. **An efficient context-free parsing al-**
686 **gorithm.** *Communications of the ACM*, 13(2):94–
687 102. 688
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi,
689 Sanchit Agarwal, Shuyang Gao, Adarsh Kumar,
690 Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020.
691 **MultiWOZ 2.1: A consolidated multi-domain dia-**
692 **logue dataset with state corrections and state track-**
693 **ing baselines.** In *Proceedings of the Twelfth Lan-*
694 *guage Resources and Evaluation Conference*, pages
695 422–428, Marseille, France. European Language
696 Resources Association. 697
- Claire Gardent, Anastasia Shimorina, Shashi Narayan,
698 and Laura Perez-Beltrachini. 2017. **The WebNLG**
699 **challenge: Generating text from RDF data.** In *Pro-*
700 *ceedings of the 10th International Conference on*
701 *Natural Language Generation, INLG 2017, Santi-*
702 *ago de Compostela, Spain, September 4-7, 2017*,
703 pages 124–133. Association for Computational Lin-
704 guistics. 705
- Paul Grice. 1975. Logic and conversation. In *Syntax*
706 *and semantics*, volume 3, pages 41–58. Academic
707 Press. 708
- Irene Langkilde and Kevin Knight. 1998. **Generation**
709 **that exploits corpus-based statistical knowledge.** In
710 *36th Annual Meeting of the Association for Compu-*
711 *tational Linguistics and 17th International Confer-*
712 *ence on Computational Linguistics, Volume 1*, pages
713 704–710, Montreal, Quebec, Canada. Association
714 for Computational Linguistics. 715
- Rémi Lebret, David Grangier, and Michael Auli. 2016.
716 **Neural text generation from structured data with**
717 **application to the biography domain.** In *Proceed-*
718 *ings of the 2016 Conference on Empirical Meth-*
719 *ods in Natural Language Processing, EMNLP 2016,*
720 *Austin, Texas, USA, November 1-4, 2016*, pages
721 1203–1213. The Association for Computational Lin-
722 guistics. 723
- Yohan Lee. 2021. **Improving end-to-end task-oriented**
724 **dialog system with a simple auxiliary task.** In
725 *Findings of the Association for Computational Lin-*
726 *guistics: EMNLP 2021*, pages 1296–1303, Punta
727 Cana, Dominican Republic. Association for Compu-
728 tational Linguistics. 729
- Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky,
730 Michel Galley, and Jianfeng Gao. 2016. **Deep rein-**
731 **forcement learning for dialogue generation.** In *Pro-*
732 *ceedings of the 2016 Conference on Empirical Meth-*
733 *ods in Natural Language Processing*, pages 1192–
734 1202, Austin, Texas. Association for Computational
735 Linguistics. 736
- Chin-Yew Lin. 2004. **ROUGE: A package for auto-**
737 **matic evaluation of summaries.** In *Text Summariza-*
738 *tion Branches Out*, pages 74–81, Barcelona, Spain.
739 Association for Computational Linguistics. 740

741	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>Proceedings of International Conference on Learning Representations</i> .	
742		
743		
744		
745	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 <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 780–799, Seattle, United States. Association for Computational Linguistics.	
746		
747		
748		
749		
750		
751		
752		
753		
754		
755	Ximing Lu, Peter West, Rowan Zellers, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. NeuroLogic decoding: (Un)supervised neural text generation with predicate logic constraints . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4288–4299, Online. Association for Computational Linguistics.	
756		
757		
758		
759		
760		
761		
762		
763		
764	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919, Online. Association for Computational Linguistics.	
765		
766		
767		
768		
769		
770	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	
771		
772		
773		
774		
775		
776		
777	Emmanouil Antonios Platanios, Adam Pauls, Subhro Roy, Yuchen Zhang, Alex Kyte, Alan Guo, Sam Thomson, Jayant Krishnamurthy, Jason Wolfe, Jacob Andreas, and Dan Klein. 2021. Value-agnostic conversational semantic parsing . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics</i> , Online. Association for Computational Linguistics.	
778		
779		
780		
781		
782		
783		
784		
785	Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pages 8689–8696.	
786		
787		
788		
789		
790		
791	Ehud Reiter. 2022. What are the problems with rule-based NLG? https://ehudreiter.com/2022/01/26/problems-with-rule-based-nlg/ .	
792		
793		
794	Subhro Roy, Sam Thomson, Tongfei Chen, Richard Shin, Adam Pauls, Jason Eisner, and Benjamin Van Durme. 2022. BenchCLAMP: A benchmark for evaluating language models on semantic parsing . <i>arXiv:2206.10668 [cs.CL]</i> .	
795		
796		
797		
798		
	Semantic Machines, Jacob Andreas, John Bufe, David Burkett, Charles Chen, Josh Clausman, Jean Crawford, Kate Crim, Jordan DeLoach, Leah Dorner, Jason Eisner, Hao Fang, Alan Guo, David Hall, Kristin Hayes, Kellie Hill, Diana Ho, Wendy Iwaszuk, Smriti Jha, Dan Klein, Jayant Krishnamurthy, Theo Lanman, Percy Liang, Christopher H. Lin, Ilya Lintsbakh, Andy McGovern, Aleksandr Nisnevich, Adam Pauls, Dmitriij Petters, Brent Read, Dan Roth, Subhro Roy, Jesse Rusak, Beth Short, Div Slomin, Ben Snyder, Stephon Striplin, Yu Su, Zachary Tellman, Sam Thomson, Andrei Vorobev, Izabela Witoszko, Jason Wolfe, Abby Wray, Yuchen Zhang, and Alexander Zotov. 2020. Task-oriented dialogue as dataflow synthesis . <i>Transactions of the Association for Computational Linguistics</i> , 8:556–571.	799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814
	Iulian V. Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2016. Building end-to-end dialogue systems using generative hierarchical neural network models . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> .	815 816 817 818 819
	Lifeng Shang, Zhengdong Lu, and Hang Li. 2015. Neural responding machine for short-text conversation . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1577–1586, Beijing, China. Association for Computational Linguistics.	820 821 822 823 824 825 826 827
	Richard Shin, Christopher H. 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 <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021</i> , pages 7699–7715. Association for Computational Linguistics.	828 829 830 831 832 833 834 835 836 837
	David Smith and Jason Eisner. 2006. Quasi-synchronous grammars: Alignment by soft projection of syntactic dependencies . In <i>Proceedings on the Workshop on Statistical Machine Translation</i> , pages 23–30, New York City. Association for Computational Linguistics.	838 839 840 841 842 843
	Alessandro Sordoni, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. 2015. A neural network approach to context-sensitive generation of conversational responses . In <i>Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 196–205, Denver, Colorado. Association for Computational Linguistics.	844 845 846 847 848 849 850 851 852 853
	Oriol Vinyals and Quoc Le. 2015. A neural conversation model . In <i>Proceedings of ICML Deep Learning Workshop</i> .	854 855 856

- 857 Marilyn A. Walker, Owen C. Rambow, and Monica Ro-
858 gati. 2002. [Training a sentence planner for spoken](#)
859 [dialogue using boosting](#). *Computer Speech & Lan-*
860 *guage*, 16(3):409–433. Spoken Language Genera-
861 tion.
- 862 Richard S. Wallace. 2009. [The anatomy of A.L.I.C.E.](#)
863 In *Parsing the Turing Test*, pages 181–210. Springer,
864 Dordrecht.
- 865 Yue Wang, Weishi Wang, Shafiq Joty, and Steven C.H.
866 Hoi. 2021. [CodeT5: Identifier-aware unified pre-](#)
867 [trained encoder-decoder models for code under-](#)
868 [standing and generation](#). In *Proceedings of the 2021*
869 *Conference on Empirical Methods in Natural Lan-*
870 *guage Processing*, pages 8696–8708, Online and
871 Punta Cana, Dominican Republic. Association for
872 Computational Linguistics.
- 873 Joseph Weizenbaum. 1966. [ELIZA – a computer pro-](#)
874 [gram for the study of natural language communica-](#)
875 [tion between man and machine](#). *Communications of*
876 *the ACM*, 9(1):36–45.
- 877 Sam Wiseman, Stuart Shieber, and Alexander Rush.
878 2017. [Challenges in data-to-document generation](#).
879 In *Proceedings of the 2017 Conference on Empiri-*
880 *cal Methods in Natural Language Processing*, pages
881 2253–2263, Copenhagen, Denmark. Association for
882 Computational Linguistics.
- 883 Tao Yu, Rui Zhang, Heyang Er, Suyi Li, Eric Xue,
884 Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze
885 Shi, Zihan Li, Youxuan Jiang, Michihiro Yasunaga,
886 Sungrok Shim, Tao Chen, Alexander Fabbri, Zifan
887 Li, Luyao Chen, Yuwen Zhang, Shreya Dixit, Vin-
888 cent Zhang, Caiming Xiong, Richard Socher, Wal-
889 ter Lasecki, and Dragomir Radev. 2019. [CoSQL: A](#)
890 [conversational text-to-SQL challenge towards cross-](#)
891 [domain natural language interfaces to databases](#). In
892 *Proceedings of the 2019 Conference on Empirical*
893 *Methods in Natural Language Processing and the*
894 *9th International Joint Conference on Natural Lan-*
895 *guage Processing (EMNLP-IJCNLP)*, pages 1962–
896 1979, Hong Kong, China. Association for Computa-
897 tional Linguistics.
- 898 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.
899 Weinberger, and Yoav Artzi. 2020a. [BERTScore:](#)
900 [Evaluating text generation with BERT](#). In *Proceed-*
901 *ings of the 8th International Conference on Learning*
902 *Representations, ICLR 2020*.
- 903 Yizhe Zhang, Siqu Sun, Michel Galley, Yen-Chun
904 Chen, Chris Brockett, Xiang Gao, Jianfeng Gao,
905 Jingjing Liu, and Bill Dolan. 2020b. [DIALOGPT:](#)
906 [Large-scale generative pre-training for conversa-](#)
907 [tional response generation](#). In *Proceedings of the*
908 *58th Annual Meeting of the Association for Compu-*
909 *tational Linguistics: System Demonstrations*, pages
910 270–278, Online. Association for Computational
911 Linguistics.

A Human Evaluation Details

A screenshot of the MTurk interface for human evaluation is shown in Fig. 4. The inter-annotator agreements for different systems are provided in Table 3. It can be observed that the gold responses receive the highest agreements on all three questions. The QCFG-constrained decoding has slightly higher agreements than the unconstrained decoding. The QCFG random sampling receives a significantly lower agreement on “Grammatical”, which is likely because this approach may produce ungrammatical responses but people may not agree on whether these are understandable.

B Model Configurations

For SMCaFlow, we fine-tune the CodeT5 model for a fixed number of epochs (=10). For MultiWOZ, we fine-tune the model for at most 10 epochs and do early stopping based on the loss on the development set. We use the AdamW optimizer (Loshchilov and Hutter, 2019) with $\beta_1 = 0.9$ and $\beta_2 = 0.999$, using a linear learning rate scheduler with an initial learning rate of 5×10^{-5} . For decoding, we always use a fixed beam size of 5.

The CodeT5-base models used in our experiments have 220 million parameters. We used machines with 32GB V100 GPUs for model fine-tuning while the decoding experiments were carried out on CPU-only machines.

For SMCaFlow experiments, the input sequence to the LM is the string representation of the computation in the lispess format followed by its the execution result is rendered as a JSON string, e.g., “Plan: (Yield (Event.start (. . .))) Result: {“type”: “DateTime”, “value”: . . . } <s>”, where the last token is a special token to separate the input and the output. For the ablative study (group 5) in §5.3, the user utterance is prefixed to the sequence, e.g., “User: When do I have thee oil change on my car scheduled for? Plan: . . . Result: . . . <s>”.

For MultiWOZ experiments, the computation is rendered as a raw JSON string which encodes the ground-truth belief state and the predicted system act. There is no execution result for these computations.

C Qualitative Analysis

We looked at 100 randomly selected examples from the experiments on SMCaFlow from §5.2,

and compare the generated responses from both unconstrained decoding and QCFG-constrained decoding with the human-annotated gold responses provided by the dataset. We summarize the differences between the generated and gold responses in Table 4, using the following categories:

Untruth The system reports incorrect information.

Omission The system fails to mention information mentioned in the gold response.

Addition The system mentions additional (correct) information that is not mentioned in the gold response.

Minor Difference The system uses a different phrasing than the gold response that nonetheless has the same information and fluency.

Disfluency The system output is disfluent.

Annotation Error The system output is acceptable but the gold annotation contains a fluency or factuality error.

For unconstrained decoding, 57 out of 100 responses differ from the gold responses, whereas for QCFG-constrained decoding, only 51 of 100 responses differ. This result is consistent with the R@1 column of Table 1 (mismatch rates of 53% and 44% respectively on the full validation set).

As expected, the most noticeable difference is the number of Untruths reported by the unconstrained system – 19%, close to the 18% rate found in the human evaluations in Table 1. We show some examples of Untruths in Fig. 5. The QCFG-constrained system produced no Untruths. Conversely, the QCFG-constrained system produces substantially more Omissions than the unconstrained system. Of the 11 omissions produced by the constrained system, 3 are identical to the unconstrained output while 7 are on inputs for which the unconstrained output produce an Untruth. In other words, our system successfully removed the 19 Untruths by the system, but in 7 of those cases, it produced a shorter (but still factually correct) input than the preferred gold annotation for that example. We also note that the gold dataset is not consistent in how much information is included in the responses – short answers like “Looks like it” in Example C from Fig. 5 are present in the gold annotations on examples similar to Example C. Furthermore, both

Instructions
Shortcuts

Instructions

In this task, you are asked to rate the quality of a virtual assistant's response to a user's request about their calendar. Please carefully read the instructions below. You are also strongly encouraged to read an example by clicking the "More Instructions" link at the end of this page.

You need to read a dialogue exchange between a user and a virtual assistant, and you are provided with all events in the user's calendar and the time when the user made the request. Then you need to answer three questions (Q1, Q2, Q3) about the quality of the virtual assistant's response. If you have feedback about the task, please enter your response in Q4.

In the section **Person Database**, you will see a table containing information about people in the organization. We only show a subset of people for conciseness.

In the section **Event Database**, you will see a table containing all events in the user's calendar. Sometimes the table can be empty, meaning there is no event in the calendar.

The section **Timestamp** provides the date and time when the user makes the request. This information is often useful for answering Q3.

The dialogue exchange is provided in the section **Dialogue**. The user is always Damon Straeter. Note sometimes the organizer of an event in the calendar may be someone other than Damon Straeter.

The section **Questions** has three required questions (Q1, Q2, Q3) about the quality of the virtual assistant's response. For Q2 and Q3, if for some reason it is impossible to judge (e.g., when the virtual assistant's response is uninterpretable), you can choose the option "Unable to decide". If you have any feedback about this task, please enter your response in Q4.

Person Database

Name	Email	Manager
Damon Straeter	dstraeter@thenextunicorn.com	David Lax
David Lax	dlax@thenextunicorn.com	Dan Schoffel

Event Database

ID	Subject	Start Time	End Time	Duration (minutes)	Show As Status	Location	Organizer	Attendees (Accepted)	Attendees (TentativelyAccepted)	Attendees (Declined)	Attendees (NotResponded)
1	The Fall of Reach	Wed Aug 30 16:00:00 2552	Wed Aug 30 16:30:00 2552	30	Busy	N/A	The Vadamee				Damon Straeter [Required]

Timestamp

Wed Aug 14 15:09:22 2019

Dialogue

Damon Straeter: Tell me who organized the fall of Reach.

Virtual Assistant: Vadamee is the organizer of "Fall of Reach".

Questions

Q1: Has the virtual assistant made any grammar errors?

No [ⓘ]

Yes but still understandable [ⓘ]

Yes and not even understandable [ⓘ]

Q2: Has the virtual assistant misunderstood the user's request?

No [ⓘ]

Yes [ⓘ]

Unable to decide [ⓘ]

Q3: Has the virtual assistant provided any incorrect information as judged using the database and timestamp?

No [ⓘ]

Yes [ⓘ]

Unable to decide [ⓘ]

Q4 (Optional): Do you have any feedback on this task?

Figure 4: A screenshot of the MTurk interface for human evaluation.

1009 systems produce more Additions than Omissions, 1026
 1010 indicating that there is not a systematic bias to- 1027
 1011 wards shorter answers overall. In future work, the 1028
 1012 model could be made to select more descriptive re- 1029
 1013 sponses by adding a brevity penalty in the decoder 1030
 1014 or by weighting the QCFG productions, so that re- 1031
 1015 sponses are scored not only by the LM but also by 1032
 1016 the QCFG. 1033

1017 D Limitations and Future Direction

1018 Authoring transduction rules is relatively easy but 1035
 1019 may be still labor intensive for complex domains. 1036
 1020 Future work might explore (semi-)automatically 1037
 1021 deriving transduction rules from data or learning 1038
 1022 to synthesize them from domain specifications, or 1039
 1023 curating a collection of domain-general transduc- 1040
 1024 tion rules which can be directly used in new do- 1041
 1025 mains.

Another direction would be to extend the dataflow transduction rules so they can encode pragmatic knowledge and context-dependent policies. For example, a dataflow transduction rule could call a neural network to assess the suitability of applying the rule to a given node in the dataflow graph, and weight the resulting QCFG production accordingly.

E Dataset License

The SMCaFlow dataset is distributed under the CC BY-SA 4.0 license. To the best of the authors knowledge, the MultiWOZ datasets were released under MIT license as shown in <https://github.com/budzianowski/multiwoz>. Our experiments follow the intended use of these datasets, which is to advance research in dialogue systems.

System	Grammatical	Relevant	Truthful
QCFG Random Sampling	.58	.75	.71
Unconstrained Decoding	.86	.71	.71
QCFG-Constrained Decoding	.90	.78	.76
Gold	.95	.81	.80

Table 3: The inter-annotator agreements for different systems.

	Unconstrained	Constrained
Untruth	19	0
Omission	3	11
Addition	17	18
Minor Difference	10	13
Disfluency	1	1
Annotation Error	7	8
Total	57	51

Table 4: Classification of differences between generated responses and human-annotated gold responses on 100 randomly sampled examples from the SMCaFlow dataset. Details are provided in Appendix C.

User Utterance	Reference Agent Utterance	QCFG-Constrained decoding	Unconstrained decoding
A. Do I have any events scheduled on Sunday the 3rd?	No, I didn't find any events on March 3rd.	No, I didn't find any events on March 3rd.	No, I didn't find any events on May 3rd .
B. Does Pikachu fly in tomorrow?	Yes, I found one event matching "fly in" with Pikachu tomorrow. It's "fly in" from 4:00 to 4:30 PM.	Yes, I found one event matching "fly in" with Pikachu tomorrow. It's from 4:00 to 4:30 PM.	Yes, I found one event matching "fly in" with Pikachu tomorrow. It's "Fly In" from 5:00 to 5:30 PM .
C. Can you find me any events during last March?	Yes, I found one event last March. It was "Dancing in Athenry" on March 17, 2019 from 4:00 to 4:30 PM.	Looks like it.	Yes, I found one event last March. It was " Erin Go Bragh " on March 17, 2019 from 5:00 to 5:30 PM.

Figure 5: Example predictions from fine-tuned CodeT5 model with QCFG-constrained decoding and with unconstrained decoding. In all the examples shown, outputs from unconstrained decoding are untruthful to the database due to content hallucination even though the model has access to the correct execution results as part of the input. We observe that in a few cases, the constrained model prefers truthful but pragmatically unhelpful omissions like such as "Looks like it" (in Example C) compared to a more specific response.