

Teaching Models new APIs: Domain Agnostic Simulators for Task Oriented Dialogue

Anonymous ACL submission

Abstract

We demonstrate that large language models are able to simulate Task Oriented Dialogues in novel domains, provided only with an API implementation and a list of goals. We show these simulations can formulate online, automatic metrics that correlate well with human evaluations. Furthermore, by filtering for dialogues where goals are met, we can use simulation to repeatedly generate training data and improve the quality of the dialogues themselves. With no human intervention or domain-specific training data, our simulations bootstrap end-to-end models which achieve a 37% error reduction over baseline in previously unseen domains. By including as few as 32 domain-specific conversations, bootstrapped models can match the performance of a fully-supervised model with $10\times$ more data.

1 Introduction

Virtual Assistants have become ubiquitous in modern life (Acharya et al., 2021). However, building these Task Oriented Dialogue (TOD) systems is laborious, requiring significant data collection and engineering resources to add support for a novel domain. As such, methods which can generalize, learn from limited examples, and require fewer engineering resources are highly desirable (Shi et al., 2019; Shah et al., 2018; Acharya et al., 2021).

To this end, works have previously identified User Simulators, wherein a model is used to emulate a human user in place of a real one, as a means of addressing these problems. User Simulators have been used to evaluate (Walker et al., 1997, 2000; Schatzmann et al., 2005) and improve Assistant models by providing additional training data (Shah et al., 2018; Acharya et al., 2021) and reward signals for Reinforcement Learning methods (Fazel-Zarandi et al., 2017; Su et al., 2018; Shi et al., 2019). Typically, these User Simulators are either limited to enhancing existing domains (Fazel-

Zarandi et al., 2017) or utilize specialized and manually engineered rules or templates for novel domains (Shah et al., 2018; Shi et al., 2019). User Simulators have often required post-hoc human intervention to ensure quality (Shah et al., 2018).

In this work, we show that modern Large Language Models (Radford et al., 2018, 2019; Lewis et al., 2020) are capable of reasonably generating dialogues when equipped with an API implementation and a desired User goal. We observe the quality of these dialogues increases with the quality of the models. Furthermore, we observe that simulation success is a strong discriminator of Assistant performance and dialogue quality.

We describe a method for bootstrapping User and Assistant models for previously unseen dialogue domains. We use Task Success, which can be automatically measured in fully synthetic dialogues, to discriminate between high- and low-quality dialogues. By adding successful dialogues back into the training set and retraining the model, we bootstrap an Assistant model without the use of any domain-specific training data, hand-engineered rules, Natural Language templates, or humans-in-the-loop. Our methodology shows improvements in both zero-shot and full-shot settings.

Furthermore, we show that we can use Task Success as a method for automatically identifying the weakest areas of our model, and employ Active Learning (Tur et al., 2003; Olsson, 2009) to enhance performance. By additionally including as few as 32 domain-specific training examples, we can match the performance of a fully-supervised baseline provided with $10\times$ more data.

2 End-to-End TOD Conversation Setup

A high-level illustration of our simulation system is shown in Figure 1. Our simulation system consists of three main components: a *User* model, an *Assistant* model, and an *API Implementation*. A conversation consists of repeated turns between

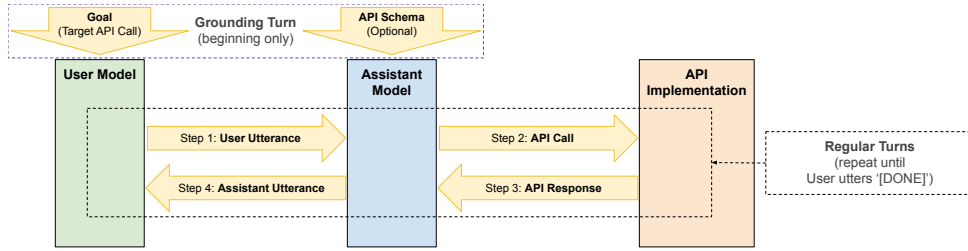


Figure 1: Illustration of our Simulation system. Initially, the User model is given a Goal API call and the Assistant is optionally provided an API Schema as grounding. Both models engage in dialogue until the User terminates the conversation. A dialogue is successful if the Assistant generates the correct API call by the end of the dialogue. In-arrows designate inputs to an entity; out-arrows designate what it generates.

081 these components. While traditional TOD systems
 082 model conversations with a combination of intent
 083 detection, belief state tracking, and policy (Jurafsky and Martin, 2009), we employ a more modern
 084 setup (Rastogi et al., 2019) where the Assistant
 085 must both make API calls at the right time and
 086 translate returned API responses into Natural Lan-
 087 guage utterances for the User. This formulation is
 088 particularly amenable to modern End-to-End (E2E)
 089 approaches based on pretrained Language Models
 090 (Ham et al., 2020; Peng et al., 2020; Hosseini-Asl
 091 et al., 2020).
 092

093 In order to guide the conversation, the User is
 094 given a *Goal* as its first turn. This Goal consists of
 095 a complete API call (e.g. intent, slot names, and
 096 slot values) serialized as a string. This is explic-
 097 itly not shown to the Assistant model to prevent
 098 information leakage. With this grounding, the User
 099 provides a natural language utterance to the Assis-
 100 tant. The Assistant optionally generates a serialized
 101 API call string and sends it to the API Implementa-
 102 tion. The API Implementation returns a serialized
 103 API response back to the Assistant based on the
 104 call, including a sentinel value for failed calls. The
 105 Assistant generates a natural language utterance to
 106 the User with this response.

107 Conversations continue in this repeated fashion
 108 of *User utterance* \rightarrow *Assistant API call* \rightarrow *API*
 109 *Implementation response* \rightarrow *Assistant utterance*
 110 until the User generates a ‘[DONE]’ token. A con-
 111 versation is said to be *successful* if the Assistant
 112 generates an API call equal to the Goal given to the
 113 User. To have a successful conversation, we expect
 114 a User to ground its generations on the Goal, the
 115 Assistant to ground on API Responses, and both of
 116 these to ground on prior utterances. We later show
 117 that simulations’ *Task Success Rate* (TSR), the suc-
 118 cess averaged over a large number of goals, is a
 119 strong proxy for the quality of dialogues generated.

120 In addition to giving Goals to the User, our sys-
 121 tem optionally allows giving the Assistant an *API*
 122 *Schema* on the first turn; we use this in Sec 5. An
 123 API Schema consists of the signature of the Goal
 124 (e.g. intent and slot names) without slot values. As
 125 we will later see, Schemas (combined with Task
 126 Success) enable us to bootstrap models in unseen
 127 domains. We take care to *not* use Schemas for eval-
 128 uations unless explicitly stated, however, since this
 129 implies a pipelined system with a precursor intent
 130 detection step to select the Schema.

3 Related Work 131

132 User Simulators have a long and successful his-
 133 tory spanning many years. They have been used
 134 widely for both evaluation and Assistant improve-
 135 ment. Formulations have varied across Rule-based,
 136 Agenda-based, and End-to-End approaches.

137 A wide range of works have explored using
 138 User Simulators as a proxy for Assistant evalua-
 139 tion (Schatzmann et al., 2005) or predicting real
 140 user satisfaction (Walker et al., 2000; Ai and Weng,
 141 2008; Jung et al., 2009; Li et al., 2016; Crook and
 142 Marin, 2017). Performance of simulators are often
 143 measured either via Task Success Rate (Gür et al.,
 144 2018; Kreyssig et al., 2018) or by cross-examining
 145 them against a wide-variety of Assistant systems
 146 (Schatzmann et al., 2005). Prior works explored
 147 the use of Language Models (LMs) as User Simula-
 148 tors, but observed that models of the time had poor
 149 adherence to goals (Georgila et al., 2006; Crook
 150 and Marin, 2017). Accordingly, most simulators
 151 have been Agenda-based, following templates and
 152 agendas to fulfill their goals (Schatzmann et al.,
 153 2007; Shah et al., 2018). In contrast, we find that
 154 modern pretrained Language Models are able to
 155 ground on and follow goals very well.

156 Engaging User Simulators and Assistant systems
 157 in synthetically generated dialogues leads to a nat-

ural reward function, and many works have used simulators in order to optimize a Reinforcement Learning policy (Schatzmann et al., 2007; Fazel-Zarandi et al., 2017; Peng et al., 2017; Su et al., 2018; Gür et al., 2018; Kreyssig et al., 2018). Such approaches are particularly used for optimizing the policy component of pipeline-based systems (Fazel-Zarandi et al., 2017), and frequently rely on the use of Natural Language Generation (NLG) templates over dialogue acts (Fazel-Zarandi et al., 2017; Shah et al., 2018; Shi et al., 2019; Kreyssig et al., 2018; Acharya et al., 2021). Our work instead utilizes fully lexicalized, E2E models for both the User and the Assistant models, without the need for agendas, dialogue acts, or NLG templates.

Most closely related to our work is that of Shah et al. (2018) and Acharya et al. (2021). Similar to our approach, both use User Simulators and Schemas to generate synthetic data and add them to a training dataset. However, Shah et al. (2018) uses an additional stage of human-annotation in order to identify failed dialogues and paraphrases, while we show that Task Success can successfully identify failures without human involvement. Acharya et al. (2021) generalizes in a few-shot setting using a similar loop as ours, but rely on templates for NLG and human paraphrases; we use sampling to achieve diversity (Holtzman et al., 2020). In contrast to prior methods that require human selection, template generation, or paraphrasing, we can both identify failures automatically and generate conversational variety automatically.

End-to-End systems for TOD have had a recent surge in interest (Asri et al., 2016; Bordes et al., 2016; Liu and Lane, 2018; Rastogi et al., 2019; Ham et al., 2020; Hosseini-Asl et al., 2020; Peng et al., 2020; Lin et al., 2020), thanks to the success of pretrained models (Devlin et al., 2019; Radford et al., 2019; Lewis et al., 2020). E2E models promise to lower the cost of annotation by replacing traditional pipeline models with text-in-text-out and adjacent external API calls (Rastogi et al., 2019; Byrne et al., 2021). Compared to these works, we leverage pretrained models for User models in order to produce simulations, rather than only modeling Assistants. We also leverage Schema-based grounding techniques (Rastogi et al., 2019; Balaraman and Magnini, 2021), which may be viewed as a form of in-context prompting methods (Brown et al., 2020; Schick and Schütze, 2020; Wang et al., 2021; Wei et al., 2021).

4 Evaluating Synthetic Dialogue

We experiment to validate the appropriateness of our Simulator setup. We generate synthetic conversations with a variety of different models and validate that traditional automated offline metrics, our TSR metric, and human evaluation all correlate and align with general expectations. Note that our objective is to verify that E2E models produce reasonable synthetic conversations using only goals, and not State-of-the-Art performance. To this end, we expect more ‘powerful’ models to outperform weaker ones, and primarily aim to validate the directional correlation across models.

4.1 Experimental Setup

Dataset We focus our efforts on the Google Schema Guided Dialogue (Google SGD) dataset (Rastogi et al., 2019). Google SGD is a large TOD dataset with emphasis on zero-shot incorporation of new skills. Models may make Service Calls (API Calls) which return responses. As part of the data, models may have access to *API Schemas*, which are call signatures of services and intents. The held-out test set contains several services and intents which are not shown in the training data. In this section, we do *not* employ the API Schemas. Instead, our Assistant models must learn the underlying Schemas directly from the data.

Models We experiment with four model architectures: LSTM (Hochreiter and Schmidhuber, 1997), LSTM with Attention (Bahdanau et al., 2014), GPT2 (Radford et al., 2019), and BART (Lewis et al., 2020) with R3F (Aghajanyan et al., 2021). We expect models later in this list to be more powerful. Our GPT2 implementation closely resembles the setup of SimpleTOD (Hosseini-Asl et al., 2020), while our BART model roughly mirrors the implementation of MinTL (Lin et al., 2020).

We fine-tune on Google SGD, using the original splits from Rastogi et al. (2019). For each model type, we fine-tune separate User and Assistant models. We generate synthetic conversations as described in Sec 2. We seed conversations with single API calls extracted from a specified fold, one call per conversation, as goals. We mock API Implementations via a lookup table with fully realized API calls as keys and corresponding API responses as values; this lookup table is populated directly from the dataset and returns a sentinel failure value if the Assistant requests any API call not in the dataset.

Evaluation Metrics We report two metrics for the Assistant: Joint Goal Accuracy (JGA)¹ and BLEU score. We additionally evaluate our simulation quality by Task Success Rate over the goals of the Valid and Test sets. In Google SGD, the Test set contains Out-of-Domain (OOD) examples that do not appear in the Train or Valid sets. As such, considering both Valid and Test gives us some estimate of OOD performance. While we aim to ensure that the models are well trained, e.g. comparable to results in the literature, our main objective to examine the correlation between Task Success measured on synthetic data and established offline metrics.

Human Evaluation We use ACUTE-Eval (Li et al., 2019) for human evaluation. ACUTE-Eval is a pairwise evaluation in which an annotator is shown two dialogues and asked a question about which they prefer. We ask annotators the question “Which Assistant would you rather use yourself?” As recommended by Li et al. (2019), we use a manually-curated control pair comparing an artificially repetitive dialogue with a gold dialogue from Google SGD; annotators who failed to identify the gold dialogue were removed. We select a random subset of 400 goals derived from the Test set of Google SGD and generate synthetic conversations using a fixed BART model for the User and the different model architectures for the Assistant. We only present User and Assistant utterance turns to annotators (hiding any API calls). We collect pairwise annotations between each Assistant model architecture described above, as well as the gold dialogues from the original dataset (Human). Annotators were presented with conversations with the same goal when comparing model-generated conversations. We measure the fraction of times each model architecture (or Human) was preferred in its pairwise match up, and compare all possible pairs of model architectures. The annotators were also asked to provide justification for their selections. An image of the annotator UI is included in Appendix A.1.

4.2 Results

Automatic and Qualitative Evaluation Results are shown in Table 1. We find that performance of all metrics improves monotonically in the direction

¹We deviate from original Google SGD evaluation here and report JGA on the *API calls* rather than the *belief state*. An example receives a JGA score of 1 iff it generates the exact API call perfectly. Note that the majority of turns do not have API calls, so majority baseline is about 0.71.

Model	User	Assistant		Simul. TSR	
	BLEU	JGA	BLEU	Valid	Test
LSTM	.058	.777	.123	.042	.042
LSTM+Attn	.078	.833	.183	.302	.169
GPT2	.093	.869	.223	.474	.307
BART	.116	.897	.252	.583	.352

Table 1: **Automatic Metrics of varied Model Architectures** tested on the original Google SGD split. TSR increases along with offline metrics, but shows greater discrimination than offline metrics.

		Win %				
		L	A	G	B	H
Lose %	LSTM		<i>.48</i>	<i>.55</i>	<i>.57</i>	<i>.80</i>
	Attn	<i>.52</i>		<i>.60</i>	<i>.63</i>	<i>.83</i>
	GPT2	<i>.45</i>	<i>.40</i>		<i>.58</i>	<i>.81</i>
	BART	<i>.43</i>	<i>.37</i>	<i>.42</i>		<i>.75</i>
	Human	<i>.20</i>	<i>.17</i>	<i>.19</i>	<i>.25</i>	

Figure 2: **Pairwise Human Evaluations of simulations** by the different models, along with Human conversations. Entries marked in blue agree with automatic metrics, while those in red italics disagree with the automatic metrics. Bold numbers indicate statistical significance ($p < .05$, binomial test).

we expect: more modern models with better pre-training and regularization do better. We find that the magnitude of improvements from better modeling is more visible using either measure of TSR compared to using more traditional offline metrics.

Reading simulations seeded by goals from the Test set, we observe that LSTMs generated few unique dialogues and generally ignored goals, preferring to replace slots and intents with ones more frequently seen in training. Though it had very low TSR, LSTM utterances often declared successful task completion regardless. On the other hand, GPT2 and BART models ground strongly on given goals. In particular, we observe these stronger models are able to generate plausible dialogues even on goals from domains present only within the Test (but not Train or Valid) sets of Google SGD. This ability to generalize is tested rigorously in Sec 5.

Human Evaluation The results of our human evaluation are shown in Figure 2. We label the scores depending on whether they agree or disagree with our expectations from automatic metrics. We find only the LSTM v. LSTM-Attn pair disagrees with our expectations, while all other pairwise evaluations agree with our expectations. We also see

that Humans are greatly preferred over all the simulations, indicating none of our simulations are at human-level performance. However, we additionally note that preference for gold data roughly decreases as the quality of the system improves.

We manually inspect training logs and reasons for human preferences. As would be expected, successfully helping the User was often a provided reason for preferring one Assistant over another; conciseness, naturalness, and brevity were also mentioned.

5 Bootstrapping Novel Domains

In the previous section, we demonstrated that User models can be adequately grounded on unseen goals to guide our simulations at generating plausible synthetic dialogues, possibly even on unseen domains. In this section, we consider whether this synthetic data can be used to *bootstrap* models on completely novel domains.

At a high level, our approach depends on generating synthetic data, filtering the synthetic dialogues using Task Success, and re-training the simulation models while incorporating the synthetic dialogue. This process may be repeated for multiple iterations to form a feedback loop.

5.1 Experimental Setup

Pretraining & Data setup Following the work of Soloist (Peng et al., 2020) and TOD-BERT (Wu et al., 2020), we pretrain a BART model on a large number of open source Task Oriented Dialogue datasets. In early experimentation, we found pretraining improved zero-shot Out-of-Domain JGA performance by about 6 absolute points, and initialize with this pretrained model for all baselines and proposed models. For brevity, full descriptions of each dataset and relevant preprocessing may be found in Appendix A.2.

To make sure that our out of domain bootstrapping of Google SGD is truly out of domain relative to pretraining, we build two custom splits for Google SGD. We analyze domains present across all of our datasets and select four holdout domains that are unique to Google SGD. We denote all conversations that use any of these holdout domains as part of the *Out-of-Domain* split and all remaining conversations to be the *In-Domain* split. Note that as Google SGD is a dataset which has multiple domains in a given conversation, some non-holdout dialogue appear as part of conversations in the Out-

Fold	No. Diag.	No. Domains
Pretraining (train)	119,677	29
In-Domain train	13,888	16
In-Domain valid	1,966	16
In-Domain test	3,132	16
Out-of-domain train	2,303	4
Out-of-domain valid	768	4
Out-of-domain test	768	4

Table 2: **Dataset statistics** Datasets used for simulator pre-training and evaluation. In-Domain and Out-of-Domain refer to new splits described in 5.1. Pretraining statistics include those of In-Domain Google SGD. The Out-of-Domain train fold is used to sample goals for bootstrapping, but not for pretraining.

of-Domain split. Only In-Domain Google SGD is used for pretraining. All domains are listed in the Appendix. Final statistics of pretraining and evaluation data are provided in Table 2.

Bootstrapping Procedure We use a *Schema-Aware* models in order to generate synthetic training data, as we find that Schema Aware models are necessary for zero-shot and few-shot domain generalization. We use goals from the Train split of Out-of-Domain Google SGD, to generating grounding for synthetic training data. In order to increase data diversity and prevent overfitting, we use Nucleus generation (Holtzman et al., 2020; $p = 0.9$). We generate 20 synthetic conversations for a given goal and retain only successful conversations. These successful conversations make up the synthetic data that we use for fine-tuning, with 10% withheld and used for model selection for early stopping.

We fine-tune both Schema-Aware (for data generation) and Schema-Agnostic (for evaluation) versions of our models on this synthetic data. Recall that Schema-Aware models have the User intent given to them, and therefore only Schema-Agnostic models may be used for evaluation. Models are fine-tuned incrementally – the best model from the previous iteration acts as initialization for the next iteration – and synthetic data is accumulated across iterations. During early experimentation, we found that fine-tuning a single model on both User and Assistant roles generally performed better than fine-tuning separate models and use this multitask setup for all of our experiments. Additionally, we find that multitasking on In-Domain data alongside synthetic data helps prevent overfitting, and include it in all experimental conditions.

Experiments We perform multiple experimental comparisons for models with synthetic data.

In our first experiment, we compare the performance of Schema-Aware and Schema-Agnostic models. This is to demonstrate that providing Schemas boosts performance and raises the Task Success Rate, enabling generation feedback loops.

In our second experiment, we consider how simulations may be used to bootstrap models in a *Zero-shot setting*. In these experiments, we only provide Out-of-Domain information via goals and completely synthetic data. To show the improvement provided by simulations, we compare primarily against the Base pretrained model, which has never seen any Out-of-Domain information. To contextualize the result, we also provide results for a Fully-Supervised model upper baseline, which was fine-tuned directly on all available Out-of-Domain data.

We evaluate on Schema-Agnostic versions of these models, and report offline metrics of Out-of-Domain JGA and Assistant BLEU-4, as well as the online metric TSR. We also report offline In-Domain JGA to ensure the use of synthetic data does not harm prior learned domains.

In our third experiment, we consider how simulations enable a form of *Active Learning*. At each iteration, we evaluate performance of the model across all conversation goals, and identify the 8 schemas with the lowest overall performance. We select 8 conversations from the Out-of-Domain training set with goals matching these schemes and add them into the next iteration’s training data (along with synthetic data). We also add 8 samples to the validation data. This may be seen as *Active Learning*, where model performance is used to guide data collection at a lower cost. As a comparison, we evaluate a baseline trained with an equal number of *random* samples; this demonstrates the performance of few-shot modeling without any simulation. To contextualize performance, we also compare to a model which uses $10\times$ more few-shot samples, and a fully supervised model.

Finally, as model selection based on validation goals can be seen as a weak form of leakage, we evaluate all of our models on the held-out Test set, which contains entirely unseen goals and conversations, as well. Here, we evaluate whether our bootstrap procedure can be used in data-rich environments by applying it to a fully-supervised model as well.

Model	In-Domain		Out-of-Dom	
	JGA	TSR	JGA	TSR
Schema-Agnostic	.878	.292	.777	.000
Schema-Aware	.960	.839	.880	.369

Table 3: **Task-pretrained BART models with and without Schemas** on validation data. Schemas help Assistant models make correct API calls and are necessary for any successful Out-of-Domain simulation.

5.2 Results

Use of Schemas Results of our first experiment are shown in Table 3. Across both In-Domain and Out-of-Domain, we see a substantial rise in performance in Schema-Aware models. This is unsurprising, as providing Schemas allows the model to bypass Intent detection. Indeed, providing Schemas has a dramatic effect on Out-of-Domain performance, boosting JGA well above baseline performance, and enabling a non-zero Task Success Rate. Such successful conversations form the basis of our synthetic data during bootstrapping, and thus critical to our methodology.

Zero-shot Results for our Zero-shot experiments are shown in Figure 3, with additional metrics provided in the Appendix A.4. Overall, we find that Out-of-Domain JGA performance goes up by a total of 8.3 absolute points, or about a 37% reduction in total errors, despite having no access to domain-specific data. However, JGA plateaus after one iteration of simulation training and further iterations provide marginal negative value.

Meanwhile, TSR continues to improve for 3 iterations before eventually plateauing at approximately the level of the Fully-Supervised model; this suggests that JGA and TSR are no longer coupled, and that our bootstrapping procedure primarily optimizes its selection criteria: Task Success. To ensure the model was not simply making random API guesses to maximize TSR, we counted the number of API calls in simulation, and found the model converged on approximately 1 call per dialogue, matching the desired distribution.

We also find that Assistant BLEU does not change significantly compared to baseline, suggesting that improvements to Task Success do not translate to improvements in Natural Language Generation. Finally, we see that In-Domain JGA remains unchanged relative to the baseline, demonstrating that the addition of synthetic data does not come at the cost of performance in prior learned domains.

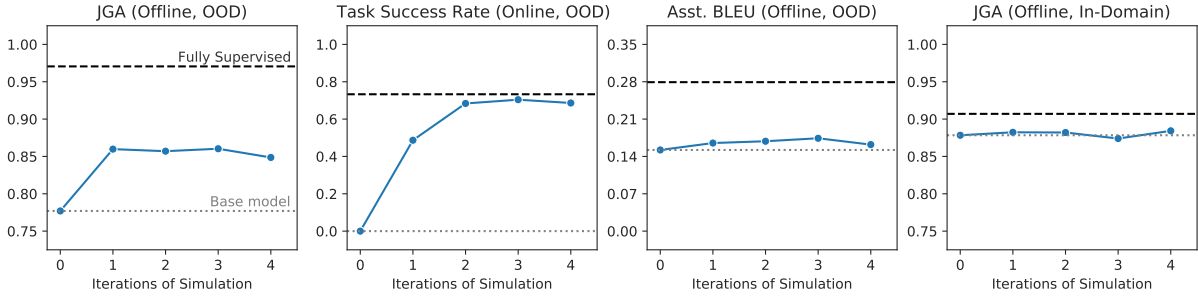


Figure 3: **Bootstrapped Validation Performance, Zero-shot.** Out-of-Domain JGA and TSR both improve through use of only synthetically generated data. BLEU and In-Domain JGA are relatively unaffected.

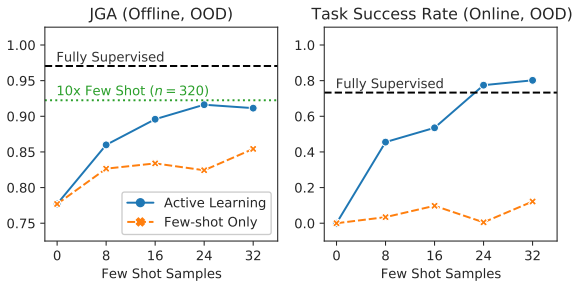


Figure 4: **Bootstrapped Validation Performance with Active Learning.** Active Learning vastly outperforms a model with an equivalent number of few-shot samples. JGA and TSR match performance of a baseline with 10× more domain-specific samples.

Active Learning Results for our Active Learning experiments are shown in Figure 4. Contrary to the Zero-shot models, we see that JGA performance consistently improves for 3 iterations, finishing with a total of 13.4 points over the Zero-shot baseline and matching the performance of Few-shot model with 320 dialogues (in green), 10× the number available to the Active Learning model. The Active Learning model also shows a large gain over the sample-equivalent Few-shot model (dashed orange), and that the gain increases with the number of samples. These results demonstrates that our use of Task Success strongly improves our sample-efficiency. TSR performance continues to improve, and eventually exceeds the fully-supervised model that was trained with 72× more data.

Although not shown, we find that Assistant BLEU, as in our Zero-shot experiments similarly does not significantly improve. This indicates that TSR is more strongly correlated with JGA, and its optimization primarily benefits NLU. In-Domain JGA also remains flat, confirming that synthetic data does not lower existing performance. Additional metrics are provided in Appendix A.4.

Model	In-Domain		Out-of-Dom	
	JGA	TSR	JGA	TSR
Base model	.829	.394	.770	.000
Simulation	.838	.454	.860	.779
Few-shot only ($n = 32$)	.835	.459	.852	.140
Active Learning ($n = 32$)	.830	.362	.911	.799
Fully Supervised	.895	.551	.973	.769
Fully Sup. + Simulation	.895	.555	.977	.847

Table 4: **Test set results.** Final performance on the held out Test-set for both In-Domain and Out-of-Domain.

Test Set Results To ensure that our methodology did not overfit via leaking of goals, we report final Test Set performance for a fully held-out set of Out-of-Domain data. For Simulation-based models, we evaluated models after 4 iterations. Results are shown in Table 4. We find that results are consistent with our earlier analysis. Zero-shot JGA improves 9 points over the baseline, and Active Learning gains 14 points over the baseline. Task Success Rate shows larger improvements, and matches the Fully-supervised baseline. Both models outperform the Few-shot only baseline. Finally, we see that our simulation procedure remains useful even in data-rich environments: adding simulations to a fully supervised model improves JGA by 0.4 absolute points (15% error reduction).

Human Evaluation We perform a final human evaluation using each the models from our experimental conditions. We repeat the ACUTE-Evals described in Section 4.2, using synthetic dialogues from each condition. Results are shown in Figure 5.

We find that all human judgements are roughly consistent with our expectations from offline evaluation, with clear wins for Active Learning over the Baseline and Few-shot models. The Human gold data outperforms all of our models, indicating further avenues for improvement. Nonetheless,

Lose %		Win %				
		B	F	Z	A	H
	<u>Base Model</u>		.60	.73	.59	.72
	Few-shot only	.40		.60	.62	.75
	<u>Zero-Shot Sim.</u>	.27	.40		.42	.66
	<u>Active Learning</u>	.41	.38	.58		.69
	<u>Human</u>	.28	.25	.34	.31	

Figure 5: **Pairwise Human Evaluations of bootstrapped models.** Entries marked in blue agree with automatic metrics, while those in red italics disagree with the automatic metrics. Bold numbers indicate statistical significance ($p < .05$, binomial test).

USER: I need a rental car.

ASST: What date and time do you need the car?

USER: I need it from the 2nd of March until the 5th.

ASST: What time do you need the car?

USER: I need it at 12:30 pm.

ASST: What city do you want to rent the car in?

USER: I want to rent a car in Fresno, CA.

ASST: I found a Dodge Charger available at 12:30 pm.

Figure 6: **Synthetic Conversation** from our final Active Learning model with *Greedy* generation. The conversation is successful, but linguistic variation is low.

the win rate of Humans decreases in our models compared to the Baseline and Few-shot models.

In analysis of annotator preferences, we find that annotator selection is generally well-correlated with TSR: Annotators preferred unsuccessful conversations over successful conversations in only about 10% of pairings. Simplicity and clarity were oftentimes given as rationale for preference in these unaligned pairings. Overall, annotators preferred Assistants that had clear communication. Many of the models learned to ask confirmation questions; this was generally liked as long as there was not too much back and forth between the User and Assistant models in doing so. Some annotators even preferred generated conversations with confirmations over the gold, human conversations. While some annotators preferred “friendlier” or “more conversational” Assistants, this was not a consistent preference; some annotators found similar conversations to be “weirdly informal,” or to “take too long to get to the point.”

6 Limitations

While our bootstrapping and Active Learning procedures do significantly improve the robustness of

models, they do surprisingly little to affect linguistic diversity, especially when using greedy generation. As a representative example, see Figure 6. In it, we observe that the conversation devolves to a simple slot-filling questionnaire, with the User beginning many utterances with “I need. . .” Manually reviewing simulated conversations found that while hallucination is very low in greedy-generated dialogues, most dialogues form roughly the slot-filling questionnaire pattern. While the synthetic conversations are plausible, their linguistic diversity is extremely low, explaining why our TSR reaches near perfect levels: the User learns to specify things as simply as possible. Furthermore, we find one of our domains (*Make payments*) has multiple instances of infinite-loops being generated, a common issue known in Neural Language Models (Holtzman et al., 2020; Welleck et al., 2019).

These examples, along with the lack of improvements in NLG metrics (BLEU scores), show that Task Success is likely to reduce the linguistic diversity of Assistants or Users, unless generation methods prone to hallucination are used. In the future, identifying automatic-filtering techniques for NLG utterances, similar to Task Success and slot filling, could help with this problem. Other generation methods, like Diverse Beam Search (Vijayakumar et al., 2016), may also be able to perform a better hallucination-diversity trade-off. Nonetheless, the offline metrics on the Test set demonstrate that our methodology does improve robustness of the NLU components of our model, as indicated by the Out-of-Domain JGA metrics.

7 Conclusion

We explored the use of pretrained Language Models as User Simulators in order to generate synthetic dialogues, and filter these models for quality using Task Success. We demonstrated our methodology can be used to improve models in zero-shot, few-shot, and full-shot manners. By incorporating Active Learning, we additionally show that our models are able to bootstrap NLU performance to that of a model with $10\times$ more training data. We encourage future work to look for improved generation methods which improve diversity without hallucination, and to find methods for automatically grading the quality of generations. Other improvements, such as the use of Schema Descriptions (Rastogi et al., 2019; Lin et al., 2021), may provide further generalization on unseen domains.

References

- 631 Anish Acharya, Suranjit Adhikari, Sanchit Agarwal,
632 Vincent Auvray, Nehal Belgamwar, Arijit Biswas,
633 Shubhra Chandra, Tagyoung Chung, Maryam Fazel-
634 Zarandi, Raefer Gabriel, Shuyang Gao, Rahul Goel,
635 Dilek Hakkani-Tur, Jan Jezabek, Abhay Jha, Jiun-
636 Yu Kao, Prakash Krishnan, Peter Ku, Anuj Goyal,
637 Chien-Wei Lin, Qing Liu, Arindam Mandal, An-
638 geliki Metallinou, Vishal Naik, Yi Pan, Shachi
639 Paul, Vittorio Perera, Abhishek Sethi, Minmin Shen,
640 Nikko Strom, and Eddie Wang. 2021. *Alexa con-
641 versations: An extensible data-driven approach for
642 building task-oriented dialogue systems*. In *Proceed-
643 ings of the 2021 Conference of the North American
644 Chapter of the Association for Computational Lin-
645 guistics: Human Language Technologies: Demon-
646 strations*, pages 125–132, Online. Association for
647 Computational Linguistics.
- 648 Armen Aghajanyan, Akshat Shrivastava, Anchit Gupta,
649 Naman Goyal, Luke Zettlemoyer, and Sonal Gupta.
650 2021. Better fine-tuning by reducing representa-
651 tional collapse. *ICLR*.
- 652 Hua Ai and Fuliang Weng. 2008. User simulation as
653 testing for spoken dialog systems. In *Proceedings
654 of the 9th SIGdial Workshop on Discourse and Dia-
655 logue*, pages 164–171.
- 656 Layla El Asri, Jing He, and Kaheer Suleman. 2016.
657 A sequence-to-sequence model for user simula-
658 tion in spoken dialogue systems. *arXiv preprint
659 arXiv:1607.00070*.
- 660 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-
661 gio. 2014. Neural machine translation by jointly
662 learning to align and translate. *arXiv preprint
663 arXiv:1409.0473*.
- 664 Vevake Balaraman and Bernardo Magnini. 2021.
665 Domain-aware dialogue state tracker for multi-
666 domain dialogue systems. *IEEE/ACM Transac-
667 tions on Audio, Speech, and Language Processing*,
668 29:866–873.
- 669 Antoine Bordes, Y-Lan Boureau, and Jason Weston.
670 2016. Learning end-to-end goal-oriented dialog.
671 *arXiv preprint arXiv:1605.07683*.
- 672 Tom B Brown, Benjamin Mann, Nick Ryder, Melanie
673 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
674 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
675 Askell, et al. 2020. Language models are few-shot
676 learners. *arXiv preprint arXiv:2005.14165*.
- 677 Bill Byrne, Karthik Krishnamoorthi, Saravanan
678 Ganesh, and Mihir Kale. 2021. *TicketTalk: To-
679 ward human-level performance with end-to-end,
680 transaction-based dialog systems*. In *Proceedings of
681 the 59th Annual Meeting of the Association for Com-
682 putational Linguistics and the 11th International
683 Joint Conference on Natural Language Processing
684 (Volume 1: Long Papers)*, pages 671–680, Online.
685 Association for Computational Linguistics.
- Paul A Crook and Alex Marin. 2017. Sequence to se-
quence modeling for user simulation in dialog sys-
tems. In *INTERSPEECH*, pages 1706–1710.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
Kristina Toutanova. 2019. *BERT: Pre-training of
deep bidirectional transformers for language under-
standing*. 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 4171–4186, Minneapolis, Minnesota. Associ-
ation for Computational Linguistics.
- Maryam Fazel-Zarandi, Shang-Wen Li, Jin Cao, Jared
Casale, Peter Henderson, David Whitney, and Al-
borz Geramifard. 2017. Learning robust dialog
policies in noisy environments. *arXiv preprint
arXiv:1712.04034*.
- Kallirroi Georgila, James Henderson, and Oliver
Lemon. 2006. User simulation for spoken dia-
logue systems: learning and evaluation. In *INTER-
SPEECH*.
- Izzeddin Gür, Dilek Hakkani-Tür, Gokhan Tür, and
Pararth Shah. 2018. User modeling for task oriented
dialogues. In *2018 IEEE Spoken Language Technol-
ogy Workshop (SLT)*, pages 900–906. IEEE.
- Donghoon Ham, Jeong-Gwan Lee, Youngsoo Jang,
and Kee-Eung Kim. 2020. *End-to-end neural
pipeline for goal-oriented dialogue systems using
GPT-2*. In *Proceedings of the 58th Annual Meet-
ing of the Association for Computational Linguis-
tics*, pages 583–592, Online. Association for Com-
putational Linguistics.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997.
Long short-term memory. *Neural computation*,
9(8):1735–1780.
- Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin
Choi. 2020. The curious case of neural text degener-
ation. *ArXiv*, abs/1904.09751.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu,
Semih Yavuz, and Richard Socher. 2020. A simple
language model for task-oriented dialogue. *arXiv
preprint arXiv:2005.00796*.
- Sangkeun Jung, Cheongjae Lee, Kyungduk Kim, Min-
woo Jeong, and Gary Geunbae Lee. 2009. *Data-
driven user simulation for automated evaluation of
spoken dialog systems*. *Comput. Speech Lang.*,
23(4):479–509.
- Dan Jurafsky and James H. Martin. 2009. *Speech and
language processing : an introduction to natural
language processing, computational linguistics, and
speech recognition*. Pearson Prentice Hall, Upper
Saddle River, N.J.
- Florian Kreyszig, Iñigo Casanueva, Paweł
Budzianowski, and Milica Gasic. 2018. Neural user
simulation for corpus-based policy optimisation of
spoken dialogue systems. *ArXiv*, abs/1805.06966.

853 Ashwin K Vijayakumar, Michael Cogswell, Ram-
854 prasath R Selvaraju, Qing Sun, Stefan Lee, David
855 Crandall, and Dhruv Batra. 2016. Diverse beam
856 search: Decoding diverse solutions from neural se-
857 quence models. *arXiv preprint arXiv:1610.02424*.

858 Marilyn Walker, Ace Kamm, and Diane Litman. 2000.
859 [Towards developing general models of usability with
860 paradise](#). *Natural Language Engineering*, 6.

861 Marilyn A. Walker, Diane J. Litman, Candace A.
862 Kamm, and Alicia Abella. 1997. [PARADISE: A
863 framework for evaluating spoken dialogue agents](#).
864 In *35th Annual Meeting of the Association for Com-
865 putational Linguistics and 8th Conference of the
866 European Chapter of the Association for Computa-
867 tional Linguistics*, pages 271–280, Madrid, Spain.
868 Association for Computational Linguistics.

869 Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao.
870 2021. [Towards zero-label language learning](#). *arXiv
871 preprint arXiv:2109.09193*.

872 Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin
873 Guu, Adams Wei Yu, Brian Lester, Nan Du, An-
874 drew M Dai, and Quoc V Le. 2021. [Finetuned lan-
875 guage models are zero-shot learners](#). *arXiv preprint
876 arXiv:2109.01652*.

877 Sean Welleck, Ilya Kulikov, Stephen Roller, Emily Di-
878 nan, Kyunghyun Cho, and Jason Weston. 2019. [Neu-
879 ral text generation with unlikelihood training](#). *arXiv
880 preprint arXiv:1908.04319*.

881 Chien-Sheng Wu, Steven C.H. Hoi, Richard Socher,
882 and Caiming Xiong. 2020. [TOD-BERT: Pre-trained
883 natural language understanding for task-oriented di-
884 alogue](#). In *Proceedings of the 2020 Conference on
885 Empirical Methods in Natural Language Processing
886 (EMNLP)*, pages 917–929, Online. Association for
887 Computational Linguistics.

A Appendix

A.1 Screenshot of Annotator UI

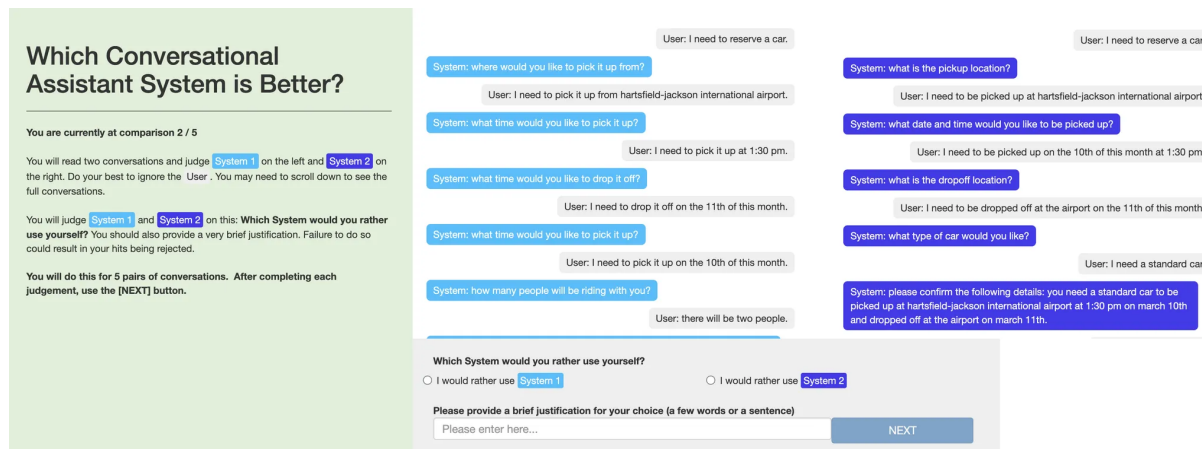


Figure 7: **Screenshot of Annotator UI.** Annotators are asked to evaluate "Which Conversational Assistant System is better" by pressing radio buttons corresponding to two presented conversations.

A.2 Pretraining Datasets

Dataset	Dial.	Dom.	Overlap with Google SGD
GoogleSGD In-Domain	18,986	16	Alarm, Banks, Buses, Calendar, Events, Flights, Hotels, Media, Movies, Music, Restaurants, Rideshare, Services, Travel, Trains, Weather
GoogleSGD Out-of-Domain	3,839	4	Home Search, Messaging, Payment, Rental Cars
MetaLWoZ	37,884	47	Banks, Buses, Events, Movies, Music, Restaurants
MSR-E2E	10,087	3	Movies, Restaurants, Taxis
MultiDoGo	19,522	6	Calendar, Flights, Media, Weather
MultiWoz	10,438	7	Attractions, Hospitals, Hotels, Restaurants, Taxis, Train
Taskmaster-1	13,215	6	Restaurants, Rideshare
Taskmaster-2	17,289	7	Flights, Hotels, Movies, Music, Restaurants
TicketTalk	23,789	1	Movies

Table 5: **Detailed statistics of datasets used in our work.** All datasets except for Google SGD Out-of-Domain used in pretraining.

We describe the different datasets that we use for pretraining. In general, we attempt to make these datasets be structured as similarly as possible to the conversations format as described in Sec 2 and generate data for separate User and Assistant models once formatted. For datasets with no API Call or Response labels, we imitate these values by accumulating dialogue state across user and assistant responses, respectively, and presenting these on appropriate turns. Other exceptions are described inline. See Table 5 for dataset statistics.

Google SGD We describe Google SGD in 4.1. We describe our method for splitting Google SGD into In-Domain and Out-of-Domain splits in 5.1.

For our Out-of-Domain split, we used Home Search, Messaging, Payment, Rental Cars as 4 holdout domains that did not have analogues in any of the other datasets. Though the "Services" and "Travel" domains do not occur explicitly in the other datasets, we do not include them in our holdout since

they include semantically similar information to the "Hospital" and "Attractions" domains of MultiWoz, respectively. For our In-Domain split, we include solely the 16 other domains of the dataset.	902
	903
As also mentioned in 5.1, since Google SGD is a dataset that contains both single-goal and multi-goal conversations, some domains of the In-Domain split are present as goals in multi-goal conversations of the Out-of-Domain split.	904
	905
	906
MetalWoz MetalWoz is a dataset constructed in a Wizard of Oz fashion across 227 tasks and 47 domains. Given a domain and a task, conversing pairs were asked to chat for 10 turns to satisfy the user’s queries.	907
	908
As this dataset does not include any annotations about API Calls, API Responses, or belief state, we pretrain on this dataset as-is and do not attempt to transform it into the format described in Sec 2. We do however split this dataset into separate User and Assistant versions.	909
	910
	911
MultiDoGo MultiDoGo is a large task-oriented dataset collected in a Wizard of Oz fashion, using both crowd and expert annotators with annotations at varying levels of granularity. We use only the data available publicly on this dataset’s open-source repository (about 20k dialogues total.)	912
	913
	914
MultiWoz MultiWoz is a dataset of single and multi-goal human-human conversations collected in a Wizard of Oz fashion. Validation and test sets contain only successful conversations while the train set include some that are incomplete. Data of the original dataset is labelled with belief states.	915
	916
	917
MSR-E2E MSR-E2E is a dataset of human-human conversations in which one human plays the role of an Agent and the other one plays the role of a User. Data is collected from Amazon Mechanical Turk.	918
	919
Taskmaster 1 Conversations in Taskmaster 1 were collected in one of two ways: spoken Wizard of Oz conversations between humans (transcribed to text) as well as written conversations from a single human in a self-dialog method. Similar to our conversations format, rather than being labelled with intents and dialog acts, conversations of this dataset are labelled with simple API arguments.	920
	921
	922
	923
Taskmaster 2 Taskmaster2 is a dataset of entirely spoken two person dialogues collected in a Wizard of Oz manner where Assistant utterances were typed by a human and then "spoken" using a text-to-speech service. Dialogues in this dataset includes those that are search and recommendations oriented, rather than purely task execution.	924
	925
	926
	927
TicketTalk TicketTalk (or Taskmaster 3) is a dataset of movie ticket dialogues collected in a self-chat manner. To induce conversational variety, crowd workers were asked to generate conversations given dozens of different instructions of different level of specificity, some purposefully including conversational errors.	928
	929
	930
	931

A.3 Hyperparameter Tables

We include hyperparameter tables for each of our models used in this paper. All models were trained using the ADAM optimizer. All models were optimized using Token Exact Match (examples with perfect greedy decoding) as an early stopping criteria.

Evaluating Synthetic Dialog Generation and its Metrics Note that as described in Sec 4, we aim to look for reasonable correlations between metrics and model architectures rather than absolute performance.

LSTM & LSTM with Attention: (1 GPU per run)

Hyperparameter	Swept Values
Learning Rate	1e{-3, -4}
Number of Layers	{1, 2, 4}
Embedding size	{256, 384}
Hidden size	{1024, 2048}
Batch size	64
Embedding Init	FastText

GPT2: (8 GPUs per run)

Hyperparameter	Swept Values
Learning Rate	1e{-5, -6}
LR Scheduler	Reduce on Plateau, Invsqrt
Model Size	124M
Text Truncate	512
Warm-up Updates	100
Batch size	4
Update Frequency	2
Gradient Clip	1

BART: (8 GPUs per run)

Hyperparameter	Swept Values
Learning Rate	1e{-5, -6}
LR Scheduler	{Reduce on Plateau, Invsqrt}
Model Size	400M
Text Truncate	512
Warm-up Updates	100
Batch size	4
Update Frequency	2
Gradient Clip	1

Bootstrapping Novel Domains Once we early stopped models on Token Exact Match, we used Task Success Rate on validation goals of Google SGD to select the best model out of a given hyperparameter sweep. However, a post hoc analysis suggests that using Token Exact Match on synthetic data fine-tunes would have worked approximately as well for the goal of improving JGA.

BART: (8 GPUs per run)

Hyperparameter	Swept Values
Learning Rate	{1e-4, 5e-5, 1e-5, 5e-6}
Model Size	400M
Batch size	4
Update frequency	8
LR Scheduler	Invsqrt
Warm-up updates	1000
Text truncate	512
Label truncate	512
Gradient Clip	0.1
Multitask Weights	1
Validation Steps	100

A.4 Additional Results

We report additional metrics on each of our models, for both offline (static, held-out data) and online (during simulation) settings.

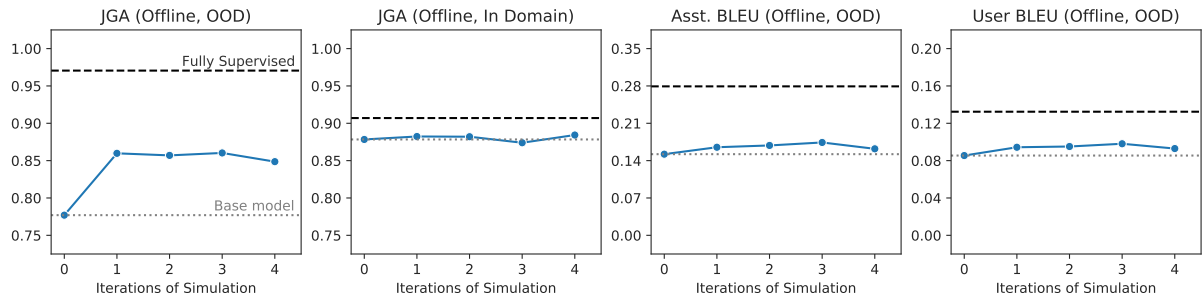


Figure 8: **Offline Bootstrapping Results.** Results of Bootstrapping on a static (held-out) offline dataset.

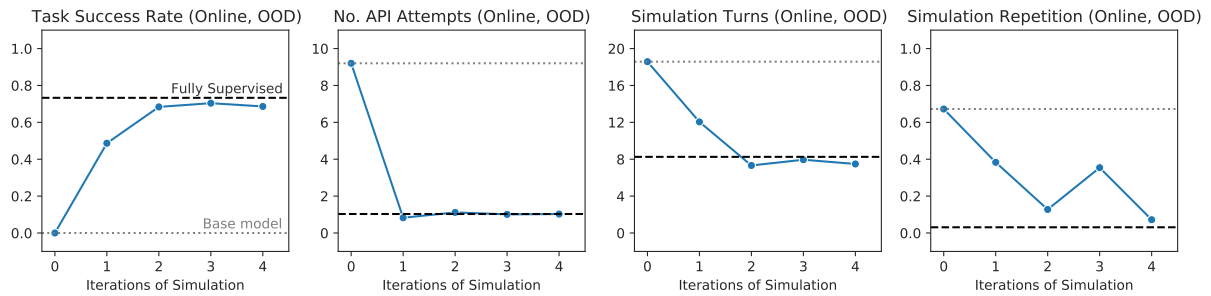


Figure 9: **Online Bootstrapping Results.** Results of Bootstrapping models on during simulations.

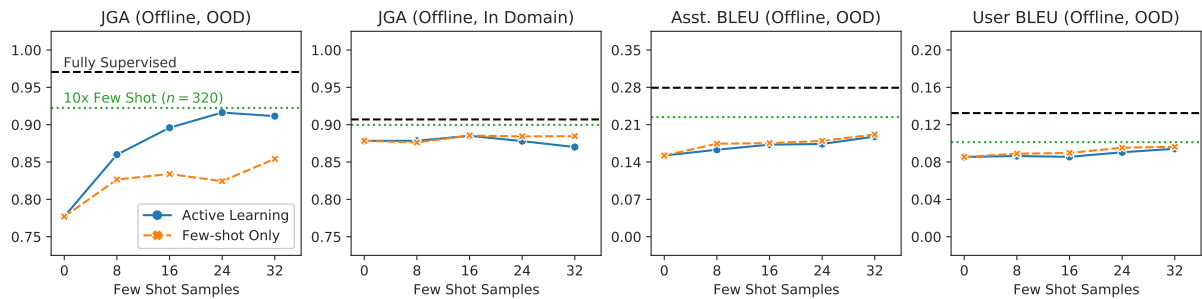


Figure 10: **Offline Active Learning Results.** Results of Active Learning on a static (held-out) offline dataset.

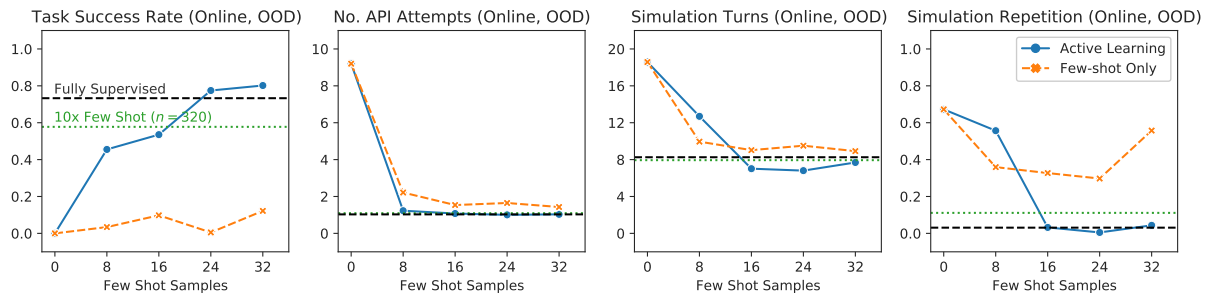


Figure 11: **Online Active Learning Results.** Results of Active Learning model during simulations.

A.5 Holdout API Analysis of Bootstrap Procedure

We take the models generated from Sec 5.1 and evaluate the models for JGA, limiting only turns to which include API calls, over each of the holdout domains. Results of this are shown in Figure 12. We observe that all holdout domains have a JGA value of zero for the Base model; this validates our selection of holdout domains. We also observe that compared to the other models, Active Learning performs much better across all domains.

Model	Domain			
	Find Homes	Payment	Rental Cars	Messaging
Base Model	.000	.000	.000	.000
Few-shot only	.603	.022	.411	.768
Zero-Shot Sim.	.876	.022	.266	.929
Active Learning	.882	.870	.623	.946

Figure 12: **JGA of individual Holdout Domains** (limited to API Call Turns only).