

BOTSTALK: Machine-Sourced Framework for Automatic Curation of Large-scale Multi-skill Dialogue Datasets

Anonymous ACL submission

Abstract

001 Previous work in open-domain chatbots has
002 introduced dialogue corpora and tasks that
003 aim to inject dialogue systems different com-
004 municative skills such as being personable,
005 knowledgeable and empathetic. With the ad-
006 vent of conversational agents grounded to
007 specific skills, a new challenge in open-
008 domain chatbots has been posed: A good
009 open-domain chatbot should retain a well-
010 rounded set of skills and seamlessly blend
011 them into a conversation. To this end, a new
012 dialogue dataset Blended Skill Talk is col-
013 lected via crowdsourcing and commonly used
014 as a benchmark for multi-skill dialogue gen-
015 eration. However, such data construction ap-
016 proach requires labor intensive manual anno-
017 tation, which severely limits their utility on
018 large-scale learning. In this work, we propose
019 BOTSTALK, a novel machine-sourced frame-
020 work, where several agents participate in a
021 conversation to automatically annotate multi-
022 skill dialogues. We then present Blended Skill
023 BotsTalk (BSBT), a large-scale multi-skill di-
024 alogue dataset of 200K conversations. Exper-
025 imental results show that our dataset can be
026 effectively used as training data for multi-
027 skill dialogue systems which require an under-
028 standing of both skill blending and grounding.
029 We also demonstrate the dataset is orthogo-
030 nally applicable to diverse learning schemes
031 such as fine-tuning and multi-task learning.

032 1 Introduction

033 A considerable progress has been made towards
034 open-domain chatbots with different desirable
035 qualities in conversation. Each of these models
036 is capable of being specialized in one communi-
037 cative skill, *i.e.*, skill grounding. A number of dis-
038 tinct large-scale datasets targeting a specific con-
039 versational skill have recently become available.
040 ConvAI2 (Zhang et al., 2018) is a dataset provided
041 for research work (Kim et al., 2020b; Majumder
042 et al., 2020a; Madotto et al., 2019) that aims to en-

dow chatbots with personas, enabling them to talk
about themselves. Wizard of Wikipedia (WoW) (Di-
nan et al., 2019) is a popular option for recent stud-
ies (Kim et al., 2020a; Zhao et al., 2020; Lian et al.,
2019) that focus on knowledgeable conversational
agents discussing topics in depth. Empathetic Dia-
logues (ED) (Rashkin et al., 2019) is also commonly
used by recent studies (Majumder et al., 2020b;
Santhanam and Shaikh, 2019) to embody empathy
in dialogue systems. Most of such skill-grounded
datasets are designed to improve a single skill, and
thus effective when models are asked to demon-
strate the targeted conversational skill.

Benefiting from the advances of these conver-
sational agents, recent research focuses on an-
other aspect of open-domain chatbots: the abil-
ity to blend various conversational skills into one
cohesive flow in a seamless manner, *i.e.*, skill blend-
ing. A good open-domain chatbot should be able
to weave multiple behaviors and skills in a single
conversation, so that it displays, for example, lis-
tening with empathy, providing knowledgeable
responses, and talking about various topics from
everyday life within a conversation (Roller et al.,
2020a; Smith et al., 2020). It should be able to adapt
to the possibilities of different users and situations
and thus use different communicative skills appro-
priately in a dialogue.

Towards this goal, there is a need to construct
a multi-skill dialogue dataset, which comprises
multi-turn dialogues that exhibit multiple skills.
While Smith et al. (2020) propose a crowdsourced
dataset Blended Skill Talk (BST) of 5K conversa-
tions as a reliable benchmark for measuring di-
alogue systems' ability at the blended objective,
it is not sufficient to build a multi-skill chatbot
due to its limited scale. Scaling up crowdsourcing
is not feasible, as it requires labor intensive man-
ual annotation and verification. Instead, automatic
curation shows promising results on large-scale
dialogue generation (Lee et al., 2021).

In this paper, we propose an automatic data curation approach that repurposes conversational agents with individual skills for generating a large-scale multi-skill dialogue dataset without additional costs or human efforts. Our main contributions are summarized as follows.

- **BOTS**TALK, a machine-sourced framework where multiple dialogue agents grounded to individual skills engage in a conversation that blends all skills together.
- **Blended Skill BotsTalk (BS^BT)**, a large-scale multi-skill dialogue dataset which contains 200K conversations blended and grounded with a number of skills derived from ConvAI2, WoW, and ED.
- Analysis and evaluation results show that our dataset can be effectively used as training resource for multi-skill dialogue systems which require an understanding of not only skill grounding but also skill blending.

2 Related Work

2.1 Skill-grounded Dialogue Datasets

Past research in open-domain chatbots has made solid strides towards dialogue systems with desirable general qualities in a conversation. Generating responses grounded to specific conversational skill has been explored in different axes, as shown in Table 1. Zhang et al. (2018) introduce ConvAI2 dataset which consists of more than 140K utterances of crowdsourced conversations to make chit-chat models more engaging and personalized by conditioning profile information on the models. Wizard of Wikipedia (Dinan et al., 2019) task aims to explore conversation informed by expert knowledge from Wikipedia and provides about 194K utterances of conversations on about 1,250 topics. Rashkin et al. (2019) constructed a dataset, Empathetic Dialogues, comprising 50K utterances of crowdworker conversations grounded in an emotional situation in order to enable a model to converse with empathy. However, it remains unclear whether models optimized for performance along specific conversational skill can retain the learned skill while blending it with other skills.

Hence, there is a clear trend in the research of open-domain chatbots, that single-skill conversation is moving to well-grounded multi-skill conversation (Smith et al., 2020; Shuster et al., 2020; Roller

et al., 2020b). In particular, Smith et al. (2020) aims to build a conversational agent who seamlessly blends being engaging and personable (Zhang et al., 2018), knowledgeable (Dinan et al., 2019), and empathetic (Rashkin et al., 2019). In order to gauge how successful a model is at this blended objective, Smith et al. (2020) collect a new multi-skill dialogue dataset of about 5K conversations, Blended Skill Talk, via crowdsourcing. While this work provides a testbed for future studies, the scale of data could hinder further progress, since training multi-skill chatbots generally requires a large-scale dataset consisting of conversations that involve multiple skills (Shah et al., 2018).

2.2 Automatic Dialogue Data Annotation

Dialogue systems research has been consistently supported by the development of new datasets (Williams et al., 2014; Mrkšić et al., 2017; Budzianowski et al., 2018). One popular approach is to collect and annotate dialogues via crowdsourcing (Zhang et al., 2018; Dinan et al., 2019; Rashkin et al., 2019; Smith et al., 2020). However, generating multi-turn dialogues in this manner requires expensive and exhausting human efforts (Shah et al., 2018; Lee et al., 2021).

Therefore, recent research seeks to facilitate open-domain chatbot development with new datasets automatically constructed by utilizing existing datasets. For example, Lee et al. (2021) create a 45K multi-modal dialogue dataset, starting with existing text-only dialogue datasets as source dialogues, and then replacing part of sentences in source dialogues with their semantically relevant images. Yang et al. (2021) propose leveraging both image-context-response triples and large scale of textual conversations for image-grounded response generation. Sun et al. (2021) propose a Human ↔ AI collaborative data collection approach for generating diverse chit-chat response to augment task-oriented dialogues and present new chit-chat based annotations to 23.8K dialogues from two popular task-oriented datasets. Kim et al. (2021b) and Vidgen et al. (2020) present a model-based dialogue collection framework and a human-and-model-in-the-loop process for generating datasets respectively.

Motivated by this line of research, in this work, we explore how large-scale multi-skill dialogue datasets can be automatically collected with minimal human efforts for data annotation.

Dataset	Dialogue episode
ConvAI2	Skill context for speaker A: I like to ski; I hate Mexican food; I like to eat cheetos; ... Skill context for speaker B: I am an artist; I have four children; I enjoy walking for exercise; ... Dialogue context A: How old are your children? B: I have four that range in age from 10 to 21. You?
Wizard of Wikipedia	Skill context for speaker A: Armadillo Skill context for speaker B: Armadillo are ... "armadillo" means "little armoured one" in ... Dialogue context A: I don't think I've ever seen an armadillo in real life! B: I've seen them at the zoo. Armadillo means little armored one in Spanish.
Empathetic Dialogues	Skill context for speaker A: My brother jump scared me while I was out playing; Terrified Skill context for speaker B: None Dialogue context A: Just got scared to death. B: Oh no. What happened?

Table 1: Example dialogues of three single-skill datasets: ConvAI2 provides each speaker persona sentences as skill contexts; Wizard of Wikipedia provides a topic and knowledge resources as skill contexts; Empathetic Dialogues provides a situation description and emotion as skill context. We only provide two turns of dialogue contexts due to the limit on the paper length.

3 Problem Formulation

In this section, we formulate the problem of multi-skill dialogue annotation and desirable characteristics for the dialogue dataset as a training resource.

3.1 Multi-skill Dialogue Annotation

Our goal is to collect a new large-scale multi-skill dialogue dataset, which can be defined as seamlessly blending various skills over the course of a multi-turn conversation. Here, inspired by Smith et al. (2020) and Sun et al. (2021), the inputs of this task are single-skill datasets, which are separately collected on a variety of skills. Let \mathbb{K} be the set of K skill types, e.g., $\mathbb{K} = \{P, K, E\}$ where P, K, E denote personality, knowledge, and empathy derived from ConvAI2, WoW, and ED, respectively. Formally, we refer to \mathcal{D}_k as a dialogue dataset with N_k dialogue episodes for skill $k \in \mathbb{K}$.

$$\mathcal{D}_k = \{(stx_{i,k}, dtx_{i,t})\}_{i=1}^{N_k} \quad (1)$$

where $stx_{i,k}$ is a skill-relevant description (i.e., skill context) for skill k and $dtx_{i,t}$ is t dialogue turns (i.e., dialogue context) derived from the skill context, as shown in Table 1.

Based on input datasets $\mathcal{D}_1, \dots, \mathcal{D}_k$, we aim to obtain a new dialogue dataset $\tilde{\mathcal{D}}$ for K skills as an output. Formally,

$$\tilde{\mathcal{D}} = \{(\tilde{stx}_i, dtx_{i,t})\}_{i=1}^{\infty} \quad (2)$$

where \tilde{stx}_i is a set of skill contexts for \mathbb{K} and $dtx_{i,t}$ is the dialogue context derived from the multiple skills. Table 3 shows a dialogue example in output dataset $\tilde{\mathcal{D}}$. We will omit the index i when dealing with a single dialogue episode.

3.2 Desirable Characteristics of Multi-skill Dialogue Datasets

By the above annotation, we aim to build a multi-skill dialogue system that uses all target skills appropriately in a conversation. For that, we lay out two criteria that a multi-skill dialogue dataset should meet as a training resource, namely **skill blending** and **skill grounding**.

Skill blending indicates that a multi-skill dataset should enable dialogue agents to exhibit different dialogue skills in a conversation flow (Smith et al., 2020; Madotto et al., 2021), while **skill grounding** emphasizes that dialogue agents should learn to maintain the same dialogue skill when appropriate (Shazeer et al., 2017). We argue that they have a trade-off relationship as it is often difficult to sufficiently represent both skill blending and grounding in a dialogue of finite length. Therefore, a desirable multi-skill conversation should be composed of short dialogue sessions specific to different skills. We note that skill grounding and blending are not contradictory, as some skill-grounded utterances imply a natural shift in skills. As an example, suppose that given an utterance “*I like sneakers because it is comfortable.*” which demonstrates skill type P, it seems reasonable to annotate an utterance with skill type K “*It is because sneakers were primarily designed for sports.*” for next dialogue turn. This example further implies that different skills can be shifted and blended naturally so that the conversational agents learn to provide reasonable responses in a multi-skill dialogue (Roller et al., 2020a).

4 BOTS TALK Framework

We now present BOTS TALK, a novel framework that automatically annotates multi-skill dialogues based on multiple single-skill dialogue datasets. The focus of our framework is to mimic a natural conversation by featuring both skill blending and grounding within a dialogue episode. Figure 1 illustrates three main phases of the framework.

4.1 Participants in BOTS TALK

In our framework, multiple participants engage in a conversation to iteratively generate the most appropriate response.

Skill Agents 🗣️ The first participants are multiple single-skill agents who annotate the appropriate skill-grounded utterances to the dialogue. Formally, based on \mathcal{D}_k for skill k , when given skill context stx_k , dialogue history dtx_t , and response space \mathbb{U} , a skill agent has dialogue models $f : (stx_k, dtx_t) \mapsto \mathbb{U}$

$$f(stx_k, dtx_t; \theta^k) \triangleq \underset{u \in \mathbb{U}}{\operatorname{argmax}} P(u|stx_k, dtx_t; \theta^k) \quad (3)$$

where θ^k is the parameters learned for skill k .

To determine response space \mathbb{U} , we design the two main functions of the dialogue agents, generator model and ranker model, parameterized as θ_{gen}^k and θ_{rnk}^k for skill k , respectively. For θ_{gen} , we aim to generate responses from response space \mathbb{U} in a token-by-token manner, and thus employ the dodecaDialogue (Shuster et al., 2020) model which is a modification of a transformer Seq2Seq architecture. On the other hand, for $\theta = \theta_{rnk}$, we consider the response space \mathbb{U} as a list of alternatives to pick the correct response, and thus employ a transformer-based retrieval architecture to score and rank response candidates in the finite set. Specifically, we use a 256-million parameter poly-encoder (Humeau et al., 2020) pre-trained on the pushshift.io Reddit dataset. Both θ_{gen} and θ_{rnk} are fine-tuned on individual single-skill datasets.

While all skill agents would simulate what response is the most appropriate conditioned on skill context set \tilde{stx} and the current dialogue context dtx_t , only one skill agent is given priority over other skill agents, to “speak” the response per dialogue turn for the dialogue annotation. We call this *active agent*. This priority may be passed to another skill agent such that the current active agent is deactivated, and another skill agent will be newly activated to speak.

Moderator Agent 😊 A critical constraint for skill agents is that neither the generator nor the ranker for skill k can learn to read other skill contexts in \tilde{stx} for different skills. For a single-skill dialogue agent, considering all possible skill contexts in a multi-skill dialogue is non-trivial. Instead, as an omniscient oracle for all skill contexts \tilde{stx} , we aim to develop another participant named moderator agent, which mediates the conversational flow for desirable skill blending and grounding. Suppose that, given an arbitrary dialogue context dtx_t , a skill agent returns a response $res_{k,t} = f(stx_k, dtx_t; \theta^k)$. Based on a set \tilde{stx} of all skill contexts and action space \mathbb{A} (i.e., approval or refusal), a moderator agent uses the decision functions $g : (\tilde{stx}, dtx_t, res_{k,t}) \mapsto \mathbb{A}$ to examine the relevance of the response with the contexts.

4.2 Phase 1: Simulate what to speak

We integrate different dialogue setups from multiple single-skill datasets as a seed information to start a conversation. Specifically, for a dialogue episode, dialogue context is initialized as an utterance pair (i.e., two-turn dialogue) via random sampling with a single-skill dataset \mathcal{D}_k , and the skill agent for skill k becomes the initial active agent. Then, for a generalizable dialogue setup, we retrieve the most relevant skill contexts from each of all input datasets $\mathcal{D}_1, \dots, \mathcal{D}_K$ by querying the seed dialogue context with a widely used IR system.¹

In the first phase of BOTS TALK, all skill agents simulate their own responses for the next dialogue turn. Formally, given a skill context set $\tilde{stx} = \{stx_1, \dots, stx_K\}$ and the current dialogue context dtx_t in a dialogue episode, a skill agent for skill k generates a plausible response $res_{k,t}$ as:

$$res_{k,t} = f(stx_k, dtx_t; \theta_{gen}^k) \quad (4)$$

where stx_k is a skill context for skill k in \tilde{stx} .

Depending on individual skills, every skill agent returns its skill-relevant response. For example, when an dtx “I love sneakers and think they are the most comfortable shoes around.” is given, the skill agent for skill P generates a response “Oh really? I like tennis shoes more than sneakers.” as $resp$, which personalizes the dialogue agent by grounding the response to a given persona. Meanwhile, the skill agents for skill K and E generate a knowledgeable response “It is because sneakers

¹We use the implementation of Chen et al. (2017).

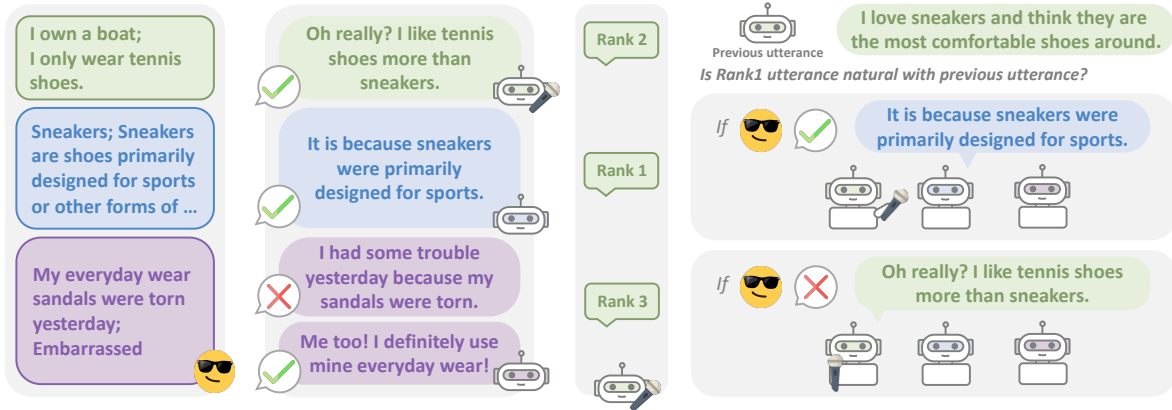


Figure 1: Illustration of BOTS TALK framework. Green, blue, and purple indicate skill types of P, K, E.

342 *were primarily designed for sports.*” as res_K and a
 343 empathetic response “*Me too! I definitely use mine*
 344 *for everyday wear!*” as res_E , respectively.

345 Note that a skill agent uses the specific skill con-
 346 text stx_k instead of \tilde{stx} for response generation.
 347 We observe that skill agents suffer from inconsis-
 348 tent and generic responses when conditioned
 349 on \tilde{stx} , as they refer less to the dialogue context
 350 dtx_t and more to skill contexts that they are not
 351 trained on. Since a skill agent aims to generate
 352 a skill-grounded response, it should take its own
 353 skill context as an input.

354 4.3 Phase 2: Check dialogue consistency

355 It is well known that neural dialogue systems lack
 356 consistency (Li et al., 2016; Welleck et al., 2019).
 357 Furthermore, as we consider different skills to-
 358 gether in a conversation, the response generated
 359 by a skill agent is more likely to be semantically
 360 in conflict with other skill contexts in \tilde{stx} . There-
 361 fore, the moderator agent, who has access to all
 362 skill contexts \tilde{stx} , is designed to maintain dialogue
 363 consistency by filtering out conflicting response
 364 candidates. A skill agent repeatedly generates new
 365 response candidates until its response $res_{k,t}$ is not
 366 contradictory to \tilde{stx} anymore.

367 Specifically, the moderator agent leverages natu-
 368 ral language inference (NLI), a task of determining
 369 whether a hypothesis sentence can be inferred
 370 from the given premise sentence. The hypothesis
 371 sentence is classified into three categories: ENTAIL-
 372 MENT (true), CONTRADICTION (false), and NEUTRAL
 373 (undetermined). A sound response $res_{k,t}$ should
 374 not be contradictory to all skill contexts \tilde{stx} . Sup-
 375 pose a stx_P is “*I wear sneakers everyday*” and a
 376 res_E is “*I had some trouble yesterday because my*
 377 *sandals were torn*”. This response is removed be-

378 cause “*yesterday because my sandals were torn*” is
 379 contradictory to “*I wear sneakers everyday*”.

380 As the moderator agent, we use a RoBERTa
 381 model (Liu et al., 2019) trained on MNLI (Williams
 382 et al., 2018)², which is widely used in fact checking
 383 systems (Kim et al., 2021a). The RoBERTa model
 384 shows 90.59% accuracy on MNLI validation set.
 385 Overall, about 50% of utterances are classified
 386 as CONTRADICTION by NLI classifier. The result
 387 demonstrates the skill agents indeed generate in-
 388 consistent responses due to the restricted access
 389 to other skill contexts. Figure 2 breaks down the
 390 result by the types of skill contexts (*i.e.* P, K, E).
 391 Out of all utterances classified as CONTRADICTION,
 392 about 70% are in conflict with other types of skill
 393 contexts. We also find that the overall proportion
 394 of utterances conflicting with stx_P is relatively
 395 high. This tendency results from the difference
 396 between skill contexts, *e.g.*, stx_P contains more
 397 distinct descriptions than stx_K and stx_E , and thus
 398 is more likely to contradict the utterance. The mod-
 399 erator agent filters out such contradictory candi-
 400 dates to preserve dialogue consistency.

401 4.4 Phase 3: Speak or pass the mic

402 Given the dialogue context and the skill contexts,
 403 the objective of the last phase is to score a set of
 404 response candidates and select a final response.
 405 To this end, we leverage the active agent and the
 406 moderator agent, taking into account a balance
 407 between skill blending and skill grounding.

408 Let \mathbb{U}_{res} be a set of response candidates
 409 $res_{1,t}, \dots, res_{k,t}$ from all skill agents. The active
 410 skill agent identifies the most appropriate response
 411 $res_t^* \in \mathbb{U}_{res}$, based on its ranker model θ_{rnk}^k , then
 412 asks to attach the selected response into the next

²Dialogue NLI (Welleck et al., 2019) is biased to ConvAI.

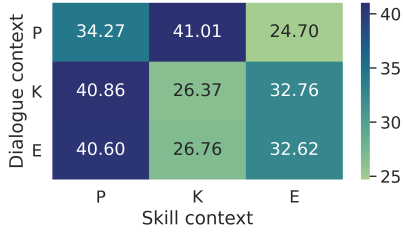


Figure 2: Percentages of utterances which are classified as CONTRADICTION via NLI classifier.

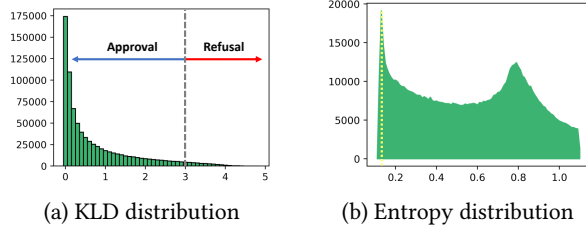


Figure 3: KL divergence between skill distributions of consecutive utterances (left) and entropy of skill distributions for all utterances (right).

dialogue context dx_{t+1} for annotation. Formally, we define such a process as:

$$res_t^* = \operatorname{argmax}_{res_t \in \mathcal{U}_{res}} P(res_t | stx_k, dx_t; \theta_{rnk}^k) \cdot g(dx_t, res_t) \quad (5)$$

where $g(dx_t, res_t) \in \{0, 1\}$ is a function of the moderator agent which determines whether res_t is approved or not.

For computing $g(dx_t, res_t)$, the moderator agent adopts a skill classifier \mathcal{P} of identifying the corresponding skill for a response. We use a BERT (Devlin et al., 2019) model trained on utterances in \mathcal{D}_k and their corresponding label k for any skill k . Once \mathcal{P} is learned, we compute $g(dx_t, res_t)$ as:

$$g(dx_t, res_t) = \begin{cases} 1, & \text{KL}(\mathcal{P}(res_{t-1}^*) || \mathcal{P}(res_t)) < \alpha \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

where res_{t-1}^* is the last utterance of dx_t and $\mathcal{P}(\cdot) \in \mathbb{R}^K$ outputs a skill distribution of an utterance/response. As the difference between two distributions (*i.e.*, KL divergence) is larger, $g(dx_t, res_t)$ is likely to have a smaller value, which is further discretized as the approval/refusal decision with a pre-defined threshold α (Figure 3a). Once the moderator agent accepts the candidate res_t^* from an inactive agent as the response, the active agent passes the mic, or the priority for annotation, to the inactive agent.

In practice, we compute entropy of all utterances based on their skill distributions in order

	Dataset	Train	Valid	Test
# Dialogues	ConvAI2	17,878	1,000	1,015
	WoW	18,430	1,948	1,933
	ED	19,533	2,770	2,547
	BST	4,819	1,009	980
	BSBT	200,000	-	-
# Utterances	ConvAI2	131,438	7,801	6,634
	WoW	166,787	17,715	17,497
	ED	40,252	5,736	5,257
	BST	41,475	8,678	8,422
	BSBT	2,000,000	-	-

Table 2: Statistics of dialogue datasets.

to investigate whether there is a room for shifting between skills. The value of entropy indicates the uncertainty of the skill type of an utterance, *e.g.*, the value of entropy becomes lower when the utterance is grounded to the specific skill and higher when the utterance is generic. As shown in Figure 3b, not all utterances are grounded to specific skill. It is evident that there are a significant amount of opportunities to shift to other skills when the entropy value of an utterance is high. This result suggests that while skill grounding, skill blending is also available in natural and reasonable way.

5 Blended Skill BotsTalk (BSBT)

5.1 Data Statistics

We collect Blended Skill BotsTalk (BSBT), a multi-skill dialogue dataset, using BOTS TALK framework. The dataset consists of 200K conversations with 2M utterances, and there are 10 utterances in each conversation. Each utterance is labeled using a skill classifier with skill annotation, which is personality from ConvAI2, knowledge from WoW or empathy from ED, including both hard label (*i.e.* skill type) and soft label (*i.e.* skill distribution). An example from BSBT is shown in Table 3.

One of the salient features of BSBT dataset is its scalability. The overall statistics of datasets used in this paper are summarized in Table 2. The BSBT dataset overwhelms other datasets with its size, since it is composed of bot-bot conversations collected through a machine-sourced approach while other datasets comprise crowdsourced human-to-human conversations. Theoretically, the size of the BSBT dataset can increase infinitely as long as seed information is provided. We ran our framework on 20 Titan RTX GPUs for 4 days to collect 200K conversations, and the data construction of 2M conversations is on progress.

Skill context from ConvAI2 I attend book club every week; My mom is my best friend; ...	Skill context from ConvAI2 I love riding roller coasters; I am a factory worker; ...
Skill context from WoW Reading process	Skill context from WoW Reading is a complex cognitive process ...
Skill context from ED I saw Stephen King on his book tour last ... & Disappointed	Skill context from ED None
Dialogue	
A: Wow that's quite a few. That's cool. (K)	
B: It's so good. It won a Pulitzer Prize. (K)	
A: Wow, what made you read it? What are you into now? (P)	
B: It was the oldest work of literature in history, because it was written in broad enough mediums. (K)	
A: Yeah, well you are probably lucky because I love reading too. (P)	
B: Yes, it is very handy to keep reading and to enjoy the parts that you like. (P)	
A: I want to read it. Hopefully I'll get a chance to see it some day. (P)	
B: I hope you will as well. Good luck! (E)	

Table 3: Sample conversation from the BS \mathbb{T} dataset. Speaker A is given five personas, one topic and a situation with an emotion (top left), while speaker B is given five personas, one topic, seven knowledge resources (top right). In the dialogue, P, K, E denotes the skill types corresponding to ConvAI2, WoW, and ED, respectively.

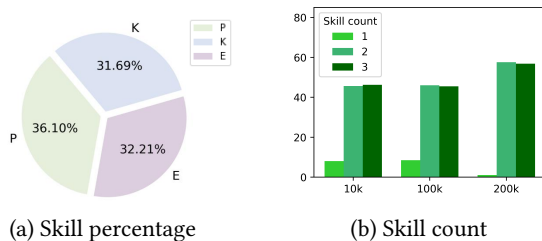


Figure 4: Illustration of skill annotation: skill percentage (left) for all utterances and the skill count per dialogue by varying the number of dialogues (right).

% annotated as:	Skill of seed utterances		
	P	K	E
P	59.79	13.76	14.89
K	18.04	72.89	10.60
E	22.15	13.33	74.50

Table 4: Percentages of utterances followed by seed utterances annotated by the skill classifier as coming from P, K, E, broken down by provenance skills of the seed utterances.

Skill Blending Figure 4 summarizes the results of skill annotation for all utterances in BS \mathbb{T} dataset. Overall, the skill annotation percentages are 36.10% for personality, 31.68% for knowledge, and 32.21% for empathy, as shown in Figure 4a. Figure 4b also shows that over 90% of the conversations demonstrate at least 2 of the 3 skills within a single conversation. This supports that the vast majority of conversations feature more than one skill, where skills of utterances are defined based on skill distribution. In addition, the tendency of skill blending in Figure 4b is stable at varying data size (10K, 100K, 200K), suggesting the efficacy of BOTS \mathbb{T} ALK on multi-skill dialogue generation.

Skill Grounding Although we focus on blending skills, the dataset should contain sufficient sessions grounded to specific skill in conversations for the model to learn the ability of skill grounding. Since the provenance skill of an utterance from original dataset is only available for seed utterances, we explore the continuity of skills based on the skill type of the utterance subsequent to seed utterances. Table 4 breaks down the results by provenance skill of the seed utterances. The fraction of utterances resembling a given dataset increases when the seed utterances are from that same dataset, and more than half of the utterances subsequent to the seed utterances are labeled the same skill type as the seed utterances.

Data Quality We perform ACUTE-Eval (Li et al., 2019), which is a popular metric for human evaluation, on BST and BS \mathbb{T} datasets. Evaluation results show that BS \mathbb{T} achieves slightly higher scores than BST, indicating the validity of BOTS \mathbb{T} ALK framework for collecting multi-skill dialogues. We provide more details in Appendix E.

5.2 Evaluation on Multi-skill Benchmark

We conduct a set of experiments to test our BS \mathbb{T} over BST benchmark. To the best of our knowledge, BST benchmark is the only multi-skill dialogue benchmark which gauges how successful a model is at blended objective. The base architecture used throughout the experiments is a 256-million parameter poly-encoder (Humeau et al., 2020) pre-trained on pushshift.io Reddit dataset. We fine-tune the base architecture on individual datasets, *i.e.*, ConvAI2, WoW, ED, BST, and BS \mathbb{T} , and consider them as our baselines.

Model	R@1	R@5	R@10	MRR
ConvAI2	75.92	94.04	97.19	83.96
WoW	67.48	89.57	94.33	77.11
ED	65.96	88.69	93.80	76.10
BST	75.92	94.76	97.83	84.14
BSBT	80.12	95.53	97.88	86.94

Table 5: Evaluation on BST benchmark.

Model	Single-skill benchmarks			Avg.
	ConvAI2	WoW	ED	
ConvAI2	88.46	79.84	47.90	72.06
WoW	57.90	90.79	45.86	64.85
ED	61.31	80.05	62.81	68.05
BST	74.13	82.12	48.11	68.12
BSBT	84.12	89.68	60.85	78.21

Table 6: Evaluation on single-skill benchmarks. Numbers in bold indicate the best performing model on the corresponding benchmark.

Model	Benchmarks			
	BST	ConvAI2	WoW	ED
MTL	78.97	86.43	90.41	59.76
MTL + BSBT50K	80.88	86.49	90.48	59.80
MTL + BSBT100K	80.94	86.71	90.63	59.92
MTL + BSBT200K	82.01	86.82	91.36	61.37

Table 7: Evaluation on benchmarks of MTL models.

The results of baselines on BST benchmark are shown in Table 5. For metrics, we measure Recall@k, or R@k, where each test example has 100 possible candidates to select from, as well as mean reciprocal rank (MRR). It is observed that multi-skill models, *i.e.*, BST and BSBT models, consistently outperform single-skill models, *i.e.*, ConvAI2, WoW, and ED models. This suggests that the single-skill models are able to do well on each of them in isolation, but struggle to seamlessly blend them over the course of a single conversation. Moreover, BSBT model outperforms all of the baselines on all automatic evaluation. This result indicates that our dataset properly works as the training resource to learn the ability to blend skills. We also provide performance of BSBT model by varying the number of dialogues on BST benchmark for scalability analysis in Appendix G.

5.3 Evaluation on Single-skill Benchmark

Table 6 summarises the results of baselines on single-skill benchmarks, *i.e.*, ConvAI2, WoW, and ED benchmarks, measured by R@1. The single-skill models each perform the best on their respective original benchmark and not as well on other benchmarks, compared to the multi-skill models,

supporting our hypothesis that single-skill agents are specialized to their corresponding skills. On the other hand, the performance of all multi-skill models is more balanced than single-skill models, in the sense that none of the single-skill models does as well averaged over the three categories (except for the ConvAI2 model doing a tiny bit better than the BST model). In particular, BSBT model performs noticeably better on all single-skill benchmarks than BST model. This suggests that BSBT is able to not only inject the ability of blending various skills but also maintain the ability for grounding specific skill.

5.4 Evaluation with Multi-task Learning

A straightforward approach of developing a multi-skill chatbot given access to multiple single-skill datasets is to multi-task on all of them during the fine-tuning step. Therefore, we consider MTL model, the poly-encoder pre-trained on pushshift.io Reddit and fine-tuned in multi-task fashion across ConvAI2, WoW, and ED. To probe the effectiveness of BSBT as training resource, we further fine-tune MTL model on BSBT datasets with different sizes (50K, 100K, 200K), respectively. Table 7 compares the performance of these models on all benchmarks, *i.e.*, ConvAI2, WoW, ED, and BST, reported by R@1 and MRR. As expected, MTL models fine-tuned on BSBT datasets with varying scales outperform MTL model for all benchmarks, indicating that BSBT is orthogonally applicable to MTL. The overall tendency also shows that the model performs better when we fine-tune the MTL model with a large scale of BSBT. Such results suggest that the scalability of BSBT is indeed crucial to model performance.

6 Conclusion

We present a novel machine-sourced approach BOTS TALK for generating multi-skill dialogues. We further propose a large-scale multi-skill dialogue dataset BSBT consisting of 200K conversations to inject a dialogue system the ability of skill blending and grounding. We demonstrate the effectiveness of our approach in comparison with several baselines by experiments on both single- and multi-skill dialogue benchmarks. Despite the inherent errors stemmed from its machine-sourced nature, our proposed data creation method can be applied when efficiently preparing datasets that cover diverse communicative skills.

References

- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gasic. 2018. Multiwoz-a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. In *Proceedings of EMNLP*.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading wikipedia to answer open-domain questions. In *Proceedings of ACL*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL*.
- Emily Dinan, Angela Fan, Adina Williams, Jack Urbanek, Douwe Kiela, and Jason Weston. 2020. Queens are powerful too: Mitigating gender bias in dialogue generation. In *Proceedings of EMNLP*.
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. Wizard of Wikipedia: Knowledge-powered conversational agents. In *Proceedings of ICLR*.
- Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux, and Jason Weston. 2020. Poly-encoders: Transformer architectures and pre-training strategies for fast and accurate multi-sentence scoring. In *Proceedings of ICLR*.
- Byeongchang Kim, Jaewoo Ahn, and Gunhee Kim. 2020a. Sequential latent knowledge selection for knowledge-grounded dialogue. In *Proceedings of ICLR*.
- Byeongchang Kim, Hyunwoo Kim, Seokhee Hong, and Gunhee Kim. 2021a. How robust are fact checking systems on colloquial claims? In *Proceedings of NAACL*.
- Hyunwoo Kim, Byeongchang Kim, and Gunhee Kim. 2020b. Will i sound like me? improving persona consistency in dialogues through pragmatic self-consciousness. In *Proceedings of EMNLP*.
- Sungdong Kim, Minsuk Chang, and Sang-Woo Lee. 2021b. Neuralwoz: Learning to collect task-oriented dialogue via model-based simulation. In *Proceedings of ACL*.
- Nyoungwoo Lee, Suwon Shin, Jaegul Choo, Ho-Jin Choi, and Sung-Hyun Myaeng. 2021. Constructing multi-modal dialogue dataset by replacing text with semantically relevant images. In *Proceedings of ACL*.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios P Spithourakis, Jianfeng Gao, and Bill Dolan. 2016. A persona-based neural conversation model. In *Proceedings of ACL*.
- Margaret Li, Stephen Roller, Ilia Kulikov, Sean Welleck, Y-Lan Boureau, Kyunghyun Cho, and Jason Weston. 2020. Don't say that! making inconsistent dialogue unlikely with unlikelihood training. In *Proceedings of ACL*.
- Margaret Li, Jason Weston, and Stephen Roller. 2019. Acute-eval: Improved dialogue evaluation with optimized questions and multi-turn comparisons. *arXiv preprint arXiv:1909.03087*.
- Rongzhong Lian, Min Xie, Fan Wang, Jinhua Peng, and Hua Wu. 2019. Learning to select knowledge for response generation in dialog systems. In *Proceedings of IJCAI*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. In *arXiv preprint*.
- Andrea Madotto, Zhaojiang Lin, Yejin Bang, and Pascale Fung. 2021. The adapter-bot: All-in-one controllable conversational model. In *Proceedings of AAAI*.
- Andrea Madotto, Zhaojiang Lin, Chien-Sheng Wu, and Pascale Fung. 2019. Personalizing dialogue agents via meta-learning. In *Proceedings of ACL*.
- Bodhisattwa Prasad Majumder, Harsh Jhamtani, Taylor Berg-Kirkpatrick, and Julian J McAuley. 2020a. Like hiking? you probably enjoy nature: Persona-grounded dialog with commonsense expansions. In *Proceedings of EMNLP*.
- Navonil Majumder, Pengfei Hong, Shanshan Peng, Jiankun Lu, Deepanway Ghosal, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2020b. Mime: Mimicking emotions for empathetic response generation. In *Proceedings of EMNLP*.
- Nikola Mrkšić, Diarmuid Ó Séaghdha, Tsung-Hsien Wen, Blaise Thomson, and Steve Young. 2017. Neural belief tracker: Data-driven dialogue state tracking. In *Proceedings of ACL*.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *Proceedings of ACL*.
- Stephen Roller, Y-Lan Boureau, Jason Weston, Antoine Bordes, Emily Dinan, Angela Fan, David Gunning, Da Ju, Margaret Li, Spencer Poff, et al. 2020a. Open-domain conversational agents: Current progress, open problems, and future directions. In *arXiv preprint*.
- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Kurt Shuster, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2020b. Recipes for building an open-domain chatbot. In *arXiv preprint*.

706	Sashank Santhanam and Samira Shaikh. 2019. Emotional neural language generation grounded in situational contexts. In <i>Proceedings of the 4th Workshop on Computational Creativity in Language Generation</i> .	A Overview	760
707			
708		In the following sections, we provide more details	761
709		on BOTS TALK framework and BS T dataset. Specifi-	762
710		cally, we lay out the details of single-skill dialogue	763
711	Pararth Shah, Dilek Hakkani-Tür, Gokhan Tür, Abhinav Rastogi, Ankur Bapna, Neha Nayak, and Larry Heck. 2018. Building a conversational agent overnight with dialogue self-play. In <i>arXiv preprint</i> .	datasets and how they are incorporated into BOT-	764
712		s TALK framework to construct BS T dataset in	765
713		Appendix B. We also provide hyperparameter set-	766
714		ups for all component models of skill agents and	767
715	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In <i>arXiv preprint</i> .	the moderator agent in Appendix C. Appendix D	768
716		presents conversation examples of BS T . In Ap-	769
717		pendix E, we provide human evaluation results of	770
718		BS T to support the validity of BOTS TALK frame-	771
719		work. In Appendix F and G, we expand the ex-	772
720	Kurt Shuster, Da Ju, Stephen Roller, Emily Dinan, Y-Lan Boureau, and J. Weston. 2020. The dialogue dodecathlon: Open-domain knowledge and image grounded conversational agents. In <i>Proceedings of ACL</i> .	perimental results on single-skill benchmark pre-	773
721		sented in Section 5.3 and conduct the scalability	774
722		analysis of BS T .	775
723			
724			
725	Eric Michael Smith, Mary Williamson, Kurt Shuster, Jason Weston, and Y-Lan Boureau. 2020. Can you put it all together: Evaluating conversational agents’ ability to blend skills. In <i>Proceedings of ACL</i> .	B Single-skill Datasets into BOTSTALK	776
726			
727		We describe the details on the single-skill dialogue	777
728		datasets used for BS T and how they are incor-	778
729	Kai Sun, Seungwhan Moon, Paul Crook, Stephen Roller, Becca Silvert, Bing Liu, Zhiguang Wang, Honglei Liu, Eunjoon Cho, and Claire Cardie. 2021. Adding chit-chats to enhance task-oriented dialogues. In <i>Proceedings of NAACL</i> .	porated into BOTS TALK framework to construct	779
730		our dataset. Example dialogues from the single-	780
731		skill dialogue datasets (<i>i.e.</i> , ConvAI2, WoW, ED)	781
732		are shown in Table 10, 11, 12.	782
733			
734	Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. 2020. Learning from the worst: Dynamically generated datasets to improve online hate detection. In <i>arXiv preprint</i> .	To integrate dialogue setups from different	783
735		single-skill datasets as a seed information, we fol-	784
736		low the basic settings for constructing a dialogue	785
737		dataset and assume a multi-turn, one-to-one con-	786
738	Sean Welleck, Jason Weston, Arthur Szlam, and Kyunghyun Cho. 2019. Dialogue natural language inference. In <i>Proceedings of ACL</i> .	versation between two speakers. We simulate turn-	787
739		taking in a conversation by switching between two	788
740		different sets of skill contexts for the input skill	789
741	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In <i>Proceedings of NAACL</i> .	context <i>stx</i> to a dialogue model <i>f</i> in skill agents.	790
742		We present details on the single-skill datasets used	791
743		to construct BS T and elaborate on how the seed	792
744		information is constructed based on them.	793
745	Jason D Williams, Matthew Henderson, Antoine Raux, Blaise Thomson, Alan Black, and Deepak Ramachandran. 2014. The dialog state tracking challenge series. <i>AI Magazine</i> , 35(4):121–124.	B.1 ConvAI2	794
746			
747		Based on PersonaChat (Zhang et al., 2018), Con-	795
748		vAI2 is a dataset of more than 140K utterances	796
749	Ze Yang, Wei Wu, Huang Hu, Can Xu, Wei Wang, and Zhoujun Li. 2021. Open domain dialogue generation with latent images. In <i>Proceedings of AAAI</i> .	from conversations in which each of paired crowd-	797
750		workers is given a role based on their persona	798
751		description and gets to know their partner. Specifi-	799
752	Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In <i>Proceedings of ACL</i> .	cally, the speaker pairs are each assigned profiles	800
753		from a set of 1155 possible personas, each consist-	801
754		ing of at least 5 profile sentences. The personas	802
755		are collected through crowdsourcing, where the	803
756	Xueliang Zhao, Wei Wu, Can Xu, Chongyang Tao, Dongyan Zhao, and Rui Yan. 2020. Knowledge-grounded dialogue generation with pre-trained language models. In <i>Proceedings of EMNLP</i> .	workers are asked to create natural, descriptive	804
757		profiles that contain typical topics of human inter-	805
758		est. Workers are also asked to keep each profile	806
759		sentence short, <i>i.e.</i> , no longer than 15 words.	807

808 Following the ConvAI2 setting, we provide two
809 different profiles composed of 5 different persona
810 sentences as the skill contexts stx_p . The personas
811 are retrieved from the ConvAI2 dataset based on
812 their relevance to the seed utterances, measured
813 by a TF-IDF retriever. The two profiles are used
814 as the input stx_p to the dialogue model f in an
815 alternating manner.

816 B.2 Wizard of Wikipedia

817 The Wizard of Wikipedia task involves discussing
818 a given topic in depth, where the goal is to both
819 engage the partner as well as display expert knowl-
820 edge (Dinan et al., 2019). The dataset consists of
821 194K utterances over 1250 topics, where each con-
822 versation begins with a randomly chosen topic. A
823 retrieval system over Wikipedia is used to retrieve
824 articles from which the dialogues are grounded
825 during the human-human crowdsourced conver-
826 sations. The topics are also crowdsourced and
827 range from commuting to Gouda Cheese to Arnold
828 Schwarzenegger. Each conversation in the dataset
829 involves two speakers named the apprentice and
830 the wizard; the apprentice aims at delving deeply
831 into a topic, while the wizard uses knowledge in
832 articles retrieved from Wikipedia to craft a rele-
833 vant reply. Specifically, given a topic derived from
834 the dialogue context, the apprentice keeps the con-
835 versation engaging and talks eagerly about a topic,
836 while the wizard responds to the Apprentice based
837 on the first paragraphs of 7 relevant Wikipedia ar-
838 ticles provided by the retrieval system.

839 In our setting, we use the simpler version of the
840 task and ignore the retrieval aspect of the task. We
841 specify the topic of the conversation as the seed
842 information of the skill agent. The skill context
843 stx_K of the Apprentice is thus defined as the given
844 topic, while stx_K of the Wizard is defined as a topic
845 and 7 relevant knowledge sources.

846 B.3 Empathetic Dialogues

847 The Empathetic Dialogues (Rashkin et al., 2019)
848 dataset includes 50K utterances of crowdworker
849 conversations grounded in an emotional situation.
850 In a conversation, one speaker describes a personal
851 situation based on an emotion label and the other
852 speaker, named the listener, displays empathy in
853 their response. Specifically, a pair of workers are
854 asked to choose an emotional word each, depict a
855 situation in 1-3 sentences based on the label, and
856 engage in a short conversation of 4-8 utterances
857 about each of the situations. Neither of the work-

858 ers, whether they be the speaker or the listener,
859 can see the emotion label and the situation descrip-
860 tion of their partner, so that they must refer only
861 to cues within the conversation for their response.

862 In our setting, we retrieve the situation descrip-
863 tion and its corresponding emotion label from the
864 Empathetic Dialogue dataset. They are used to de-
865 fine the skill context stx_E of the speaker in an ED
866 setting. Note that we do not define stx_E of the
867 listener for our framework, so that the dialogue
868 system is trained to show empathy based solely
869 on the conversation.

870 C Implementation Details

871 Our implementation is based on the ParlAI
872 toolkit,³ which is specialized in training and evalu-
873 ating dialogue systems. We will release our agents
874 and dataset for public use.

875 C.1 Skill Agent

876 In our framework, a skill agent leverages both
877 generator model and ranker model.

878 Given a stx_k and dtx as input, a generator
879 model of skill agent generates a response for the
880 next dialogue utterance. For the generator model,
881 we employ a dodecaDialogue (Shuster et al., 2020).
882 The dodecaDialogue model is a modification of
883 transformer seq2seq architecture, which has a 8-
884 layer encoder, 8-layer decoder with 512 dimen-
885 sional embeddings and 16 attention heads. We
886 fine-tune the dodecaDialogue models on ConvAI2,
887 WoW, and ED, respectively. For generative mod-
888 els, at inference time, one must choose a decoding
889 method to generate a response to the dialogue con-
890 text. In this work, we use nucleus sampling as a
891 decoding strategy.

892 Given a stx_k and dtx as input, a ranker model
893 of skill agent selects the next dialogue utterance
894 by scoring a large set of candidate responses and
895 outputting the one with the highest score. For the
896 ranker model, we employ the poly-encoder archi-
897 tecture of Humeau et al. (2020). Poly-encoders en-
898 code global features of the context using multiple
899 representations, which are attended to by each
900 possible candidate response. This final attention
901 mechanism gives improved performance over a
902 single global vector representation whilst still be-
903 ing tractable to compute compared to simply con-
904 catenating input and output as input to a Trans-
905 former. The poly-encoder has state-of-the-art per-

³<https://github.com/facebookresearch/ParlAI>

Model	R@1	R@5	R@10	MRR
Evaluation on ConvAI2 benchmark				
ConvAI2 (Zhang et al., 2018)	88.46	98.92	99.71	93.03
WoW (Dinan et al., 2019)	57.90	86.85	95.80	70.59
ED (Rashkin et al., 2019)	61.31	89.44	96.69	73.53
BST (Smith et al., 2020)	74.13	95.64	98.80	83.37
BSBT (Ours)	84.12	97.30	99.22	89.91
Evaluation on WoW benchmark				
ConvAI2 (Zhang et al., 2018)	79.84	96.97	98.84	87.62
WoW (Dinan et al., 2019)	90.79	99.28	99.66	94.67
ED (Rashkin et al., 2019)	80.05	96.25	98.37	87.34
BST (Smith et al., 2020)	82.12	97.57	98.99	89.11
BSBT (Ours)	89.68	99.38	99.77	94.09
Evaluation on ED benchmark				
ConvAI2 (Zhang et al., 2018)	47.90	76.14	85.87	60.60
WoW (Dinan et al., 2019)	45.86	74.79	85.15	58.94
ED (Rashkin et al., 2019)	62.81	88.91	94.58	74.18
BST (Smith et al., 2020)	48.11	77.09	86.96	61.04
BSBT (Ours)	60.85	87.74	94.09	72.67

Table 8: Evaluation on single-skill benchmarks, *i.e.*, ConvAI2, WoW, and ED benchmarks.

	BST (Win %)	BSBT (Win %)
Engagingness	43	57
Interestingness	47	53
Humanness	44	56

Table 9: Human evaluation for pairwise comparison between BST and BSBT datasets.

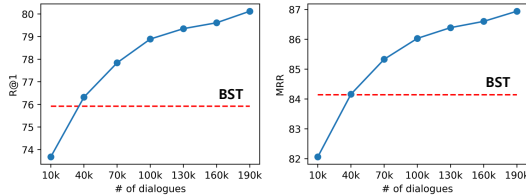


Figure 5: The effect on performance by varying the number of dialogues, reported by R@1 (left) and MRR (right).

906 performance on a number of dialogue tasks when
 907 compared to other retrieval models, and also gives
 908 comparable performance to the winning generative
 909 models on the ConvAI2 competition task in
 910 terms of human evaluation.

911 More specifically, we consider a 256M param-
 912 eter poly-encoder model. The poly-encoder is
 913 based on large pre-trained transformer models
 914 with the same architecture and dimension as BERT-
 915 base (Devlin et al., 2019), which has 12 layers,
 916 12 attention heads, and a hidden size of 768. We
 917 pre-train our poly-encoder on pushshift.io Red-
 918 dit dataset and then fine-tune on ConvAI2, WoW,
 919 and ED, respectively. We use a large number of
 920 negatives by considering the other batch elements
 921 as negative training samples, avoiding recompu-

922 tation of their embeddings. We use the Adamax
 923 optimizer without weight decay, a learning rate of
 924 5e-5 with batch size 128, epoch 8. The learning rate
 925 decays by a factor of 0.4 upon plateau of the loss
 926 evaluated on the valid set every half epoch. The
 927 best parameters are chosen based on R@1 score.

928 C.2 Moderator Agent

929 In our framework, the moderator agent leverages
 930 NLI classifier and skill classifier.

931 Given a response $res_{k,t}$ from a skill agent of skill
 932 k and the set of skill contexts \tilde{stx} , the NLI classi-
 933 fier is designed to determine whether a response
 934 candidate contradicts any of the skill contexts. For
 935 NLI classifier, we employ the public HuggingFace
 936 implementation⁴ of a RoBERTa-large model (Liu
 937 et al., 2019) fine-tuned on the Multi-Genre NLI
 938 dataset (Williams et al., 2018). The RoBERTa model
 939 shows 90.59% accuracy on MNLI validation set. We
 940 regard each response candidate res_k as hypothe-
 941 sis sentence and each skill context $stx_k \in \tilde{stx}$ as
 942 a premise sentence, then conduct unidirectional
 943 NLI between stx_k and res_k , determining whether
 944 a hypothesis sentence res_k can be inferred from
 945 the given premise sentence stx_k for all response
 946 candidates.

947 Given a response res_t , the skill classifier iden-
 948 tifies the skill of the response among all skills
 949 represented in the skill context set \tilde{stx} . For skill
 950 classifier, we employ a BERT-base (Devlin et al.,
 951 2019) model. We trained the model on utterances
 952 from ConvAI2, WoW, ED train sets and their corre-
 953 sponding skill k as labels. The model was trained

⁴<https://github.com/huggingface>

954 with a batch size of 16, a learning rate of $2e-5$ and
955 epoch 3. The BERT model shows 81.95% accuracy
956 on utterances from ConvAI2, WoW, ED test sets.

957 D Dataset Examples

958 We present a number of dialogue examples ran-
959 domly sampled from BS \mathbb{B} T in Table 13, 14, 15, 16.
960 In each dialogue episode in BS \mathbb{B} T, one speaker is
961 given five personas as stx_P , one topic as stx_K , and
962 a situation description and emotion as stx_E , while
963 another speaker is given five personas as stx_P ,
964 the topic and seven knowledge resources as stx_K ,
965 and nothing for stx_E . Each speaker is conditioned
966 on their corresponding set of skill contexts, and
967 annotates the response turn by turn.

968 E Human Evaluation

969 We perform human evaluation on BS \mathbb{B} T dataset
970 to validate our BOTS \mathbb{T} ALK framework. To this end,
971 we employ ACUTE-Eval (Li et al., 2019), which is
972 a popular metric for multi-turn dialogue evalua-
973 tion (Dinan et al., 2020; Li et al., 2020). We ran-
974 domly sample 100 dialogues from the BST and
975 BS \mathbb{B} T datasets, respectively. We then ask human
976 evaluators to compare each pair of dialogues from
977 BST and BS \mathbb{B} T datasets over three axes: engaging-
978 ness, interestingness and humanness. We provide
979 the evaluators with three questions to assess the
980 quality of the dialogues:

- 981 • **Engaging:** Who would you prefer to talk to?
982 Which version is more likely to hold your
983 attention and make you want to hear more?
- 984 • **Interesting:** Who would you say is more in-
985 teresting? Which version arouses your curios-
986 ity or tells you something new or useful?
- 987 • **Humanlike:** Who would you say sounds
988 more human? Which version is more natural
989 and personable?

990 The evaluation results are shown in Table 9. For
991 all ACUTE-Eval metrics, BS \mathbb{B} T dataset achieves
992 slightly higher win percentages over BST dataset,
993 although the difference between datasets is not sta-
994 tistically significant. Thus, our machine-sourced
995 approach BOTS \mathbb{T} ALK can be an useful alternative
996 to crowdsourcing when collecting multi-skill dia-
997 logues.

F Additional Evaluation on Single-skill Benchmarks

We provide the experimental results on single-skill
benchmarks, *i.e.*, ConvAI2, WoW, and ED bench-
marks. We report R@1, R@5, R@10, and MRR for
metrics in Table 8, which elaborates the results in
Table 6.

G Scalability Analysis

To gain more insights into the scalability of our
approach, we construct datasets at varying scales
(10K, 40K, 70K, 100K, 130K, 160K, 190K) and fine-
tune the base architecture on each of them. We
evaluate these models on BST benchmark to in-
vestigate the impact of BS \mathbb{B} T size to the model
performance. Figure 5 illustrates the performance
of BS \mathbb{B} T model in terms of R@1 when the train-
ing size of BS \mathbb{B} T dataset gradually increases. The
model fine-tuned on BS \mathbb{B} T190K dataset achieves
the best performance. The significant performance
boost from the 10K to 190K models reaffirms the
importance of large-scale training. Considering
our data curation process does not involve human
intervention (*i.e.*, no manual annotation and ver-
ification), this is a promising result showing the
potential of our dataset for large-scale supervised
learning: one can obtain datasets of arbitrary sizes
and develop multi-skill dialogue systems by lever-
aging high portion of multi-skill conversations
provided in datasets.

Example 1	
Persona 1 I like to ski My wife does not like me anymore I have went to Mexico 4 times this year I hate Mexican food I like to eat cheetos	Persona 2 I am an artist I have four children I recently got a cat I enjoy walking for exercise I love watching Game of Thrones
PERSON1: Hi PERSON2: Hello ! How are you today? PERSON1: I am good thank you, how are you. PERSON2: Great, thanks ! My children and I were just about to watch Game of Thrones. PERSON1: Nice ! How old are your children? PERSON2: I have four that range in age from 10 to 21. You? PERSON1: I do not have children at the moment. PERSON2: That just means you get to keep all the popcorn for yourself. PERSON1: And Cheetos at the moment! PERSON2: Good choice. Do you watch Game of Thrones? PERSON1: No, I do not have much time for TV. PERSON2: I usually spend my time painting; but, I love the show.	
Example 2	
Persona 1 I enjoy going to museums I am married I live in New york city I have ten nieces and nephews I am a writer	Persona 2 I like to party My major is business I am in college I love the beach i work part time at a pizza restaurant
PERSON1: Hi. How are you doing today? PERSON2: Hi, I am great just finishing up some homework. How are you? PERSON1: I am alright. I just got done writing. PERSON2: Do you write for a living or as a hobby? PERSON1: It is my living. I like culture. PERSON2: That sounds like a fun job. I am a business major but have a part time job. PERSON1: What are you going to school for? PERSON2: I am trying to get my ba in finance. PERSON1: Do you own your own company? PERSON2: No still in school work at pizza hut part time. PERSON1: Do you get free pizza? PERSON2: No but I do get a good deal.	
Example 3	
Persona 1 I work at McDonald's My parents did not want me I was raised by my uncle and aunt I am considering going to college I am looking for something different this year	Persona 2 I want to be successful I have got two more years in college I have no siblings I am a student I study law
PERSON1: My mom had me in McDonald's bathroom when she was 12. PERSON2: That is really sad. Do you have any siblings? PERSON1: No. My uncle is the manager and took me home. PERSON2: That was nice of him. I am almost done with college. PERSON1: He gave me a job at McDonalds. I was considering college but I am good now. PERSON2: I am almost done. I only have two years left in law school. PERSON1: You be smart. I would like to do something different but I team leader now. PERSON2: That sounds important. Maybe one day you will go back to school. PERSON1: I dropped out to work. Maybe you can help me sure my mom for abandoning me? PERSON2: I do not know if we would be able to do anything because is has been so long. PERSON1: Oh I am only 15 right now. Is that too old? PERSON2: Do not you have to be older than 15 to work at McDonald's? That sounds illegal.	

Table 10: 3 sample conversations from ConvAI2 dataset. In dialogue episodes, each speaker is given five persona sentences respectively.

Example 1	
Topic:	Ice cream
Wizard:	I just love ice cream. I love the types with fruits and flavours. Do you like ice cream?
Apprentice:	I love Ice cream as much as any one. I especially like Gelato, foreign ice cream!
Knowledge:	Ice cream is a sweetened frozen food typically eaten as a snack or dessert ... It is usually made from dairy products, such as milk and cream, and Bacon ice cream (or bacon-and-egg ice cream) is an ice cream generally created by adding bacon to egg custard and freezing the mixture.
Wizard:	Me too. There are some strange combinations though, have you heard of bacon ice cream? where they add bacon and even egg custard to the freezing mixture!
Apprentice:	Surprisingly bacon ice cream doesn't surprise me. That doesn't sound appealing to me, but perhaps it could be delicious ...
Example 2	
Topic:	Armadillo
Wizard:	I love animals and think armadillos are awesome with their leathery shell.
Apprentice:	I don't think I've ever seen an armadillo in real life!
Wizard:	I've seen them at the zoo. Armadillo means little armored one in Spanish.
Apprentice:	Are they native to a Spanish-speaking part of the world?
Knowledge:	Armadillos are New World placental mammals in the order Cingulata ... The word "armadillo" means "little armoured one" in Spanish. It is usually made from dairy products, such as milk and cream, and The nine-banded armadillo ("Dasypus novemcinctus"), or the nine-banded, long-nosed armadillo, is a medium-sized mammal found in North, Central, and South America.
Wizard:	Yes, they are most commonly found in North, Central, and South America
Example 3	
Topic:	Lifeguard
Apprentice:	So I am a lifeguard. Know anything about saving lives in water?
Wizard:	I'm impressed! It's a big responsibility to supervise other people's safety in the water! Tell me more.
Apprentice:	Well, I help make sure people do not drown or get injured while in or near the water!
Knowledge:	A lifeguard is a rescuer who supervises the safety and rescue of swimmers, surfers, ... Lifeguards are strong swimmers and trained in CPR/AED first aid, certified in water In some areas, the lifeguard service also carries out mountain rescues, or may function as the primary EMS provider.
Wizard:	I've heard that in some places, lifeguards also help with other sorts of emergencies, like mountain rescues! Is that part of your job too?
Apprentice:	I have! I feel like you know much about this! What brings you to know so much?
Wizard:	Oh, that's about the extent of my knowledge. I've just been around beaches and I've always admired lifeguards. I'm not a super strong swimmer myself.

Table 11: 3 sample conversations from Wizard of Wikipedia dataset. In each dialogue episode, apprentice is given a topic, while wizard is given the same topic and access to an information retrieval system over Wikipedia. For each utterance, knowledge retrieval is performed based on dialogue history, giving about 61 knowledge candidates per turn.

Example 1	
Emotion:	Terrified (label)
Situation:	My brother jump scared me while I was out playing. It was crazy bad.
Speaker:	Just got scared to death.
Listener:	Oh no. What happened?
Speaker:	My brother jumped scared me.
Listener:	lol is he younger or older?
Example 2	
Emotion:	Proud (label)
Situation:	My little dog learned to sit!
Speaker:	I finally tough my new little puppy his first trick!
Listener:	What trick did you teach him?
Speaker:	I tought him to sit for a treat, it's so cute.
Listener:	That is good, do you plan to teach him more tricks?
Example 3	
Emotion:	Apprehensive (label)
Situation:	I have to call my landlord about being late on the rent. I really don't want to have this conversation.
Speaker:	I have to make a dreadful phone call tomorrow
Listener:	Oh no, about what?
Speaker:	I'm late on my rent and I need another week. I don't want to because my landlord isnt very nice
Listener:	Oh no, I've been there done that too many times. Speaker: I don't want her to make a big deal
Example 4	
Emotion:	Content (label)
Situation:	Eating my favorite meal makes me happy.
Speaker:	I am at my best when I have my favorite meal.
Listener:	Nice.
Speaker:	I love enchiladas.
Listener:	Really?
Example 5	
Emotion:	Joyful (label)
Situation:	I have had a great week!
Speaker:	I have had a great start to my week!
Listener:	That's great. Do you think the rest of the week will be as great?
Speaker:	I hope so! It looks promising!!
Listener:	Lucky you. Are you always a positive person or it's just been an amazing week really?
Example 6	
Emotion:	Proud (label)
Situation:	I was proud when my brother finished college. He worked so hard at it.
Speaker:	I was proud of my brother when he finished school. He worked so hard at it.
Listener:	Nice, tell him congrats. What did he major in?
Speaker:	It was English.
Listener:	He should become an English teacher.

Table 12: 6 sample conversations from Empathetic Dialogues dataset. In each dialogue episode, speaker is given a situation with an emotion, while listener is given nothing as input.

Example 1	
Skill context from ConvAI2 I like to go for runs; I live by myself; ... Skill context from WoW Yoga Skill context from ED I felt really guilty when I had to leave my dog ...; Guilty	Skill context from ConvAI2 I like to listen to music; I have a boxer dog; ... Skill context from WoW Yoga; Yoga is a group of practices or disciplines ... Skill context from ED None
Dialogue	
A: Hey, do you have any kids? B: Hi! no, i shop in my spare time. do you? A: Yes, I do. Not the only shopper here. B: Nice! Do you like to swim? I am all about that exercise. A: I love swimming. I wish I could swim more though. B: I wish I could swim too. I am also much older. A: I'm not as young as I used to be. I lost weight then. I am wondering how old you are. B: I am 50 years old. Are you in school? A: Yep. I think its time to start. You don't like to wait long. B: Definitely. That is the life stage for me. I hope I can get back to age.	
Example 2	
Skill context from ConvAI2 I'm going on a cruise next month.; I love to cook.; ... Skill context from WoW Cooking Skill context from ED I tried a new recipe the other day. It did not look ...; Apprehensive	Skill context from ConvAI2 I like going barefoot.; I can knit.; ... Skill context from WoW Cooking; Cooking is the art of preparing food for ... Skill context from ED None
Dialogue	
A: Did you do any sight seeing when you last went to Dallas? B: Not really, we were there for a friends wedding. Mind you, we did go to some great restaurants. With a population of over a million. People the city has a huge number of quality places to go! A: Cool. I would have thought the same, but then again we have so many great dishes to eat. Did you enjoy any of them? B: I enjoyed them all! We stayed for more than one night, so I know how to do some great things! A: I am excited because I'm going on a cruise. B: I hope you have a great time! Cruise is such a great destination! A: There are so many different places to go on a cruise. Did you visit any other restaurants near where you were at? B: There were a lot! I saw a restaurant on my cruise tour that I had been to before. A: Was it expensive? Did you have a special restaurant? B: No, it wasn't expensive at all. We had hot dogs and rice, which I enjoyed most of the time.	

Table 13: Dialogue examples randomly sampled from BSBT dataset.

Example 3	
Skill context from ConvAI2 I go to high school.; I play ice hockey; ... Skill context from WoW Laziness Skill context from ED I participated in my high school runners' marathon ...; proud	Skill context from ConvAI2 I have a wife and two kids; I am a vegetarian; ... Skill context from WoW Laziness; Laziness is disinclination to activity or ... Skill context from ED None
Dialogue A: One job would be better then my three. B: What do you like to eat? I'm vegetarian. A: I love beef broth, grilled cheese and eggs. B: Yum! My wife and 2 kids do the same. A: Family is also my life style. Maybe we should have a child! B: But what if you didn't want a child? A: If I didn't want a child I probably wouldn't have them. B: That's true! Our kids would just be us. A: That's what I thought, but then realized I didn't want a kid either. B: I understand that, it is hard to conceive.	
Example 4	
Skill context from ConvAI2 I like cupcakes; I like to read; I work at a grocery store; ... Skill context from WoW Pecan pie Skill context from ED I really like this girl at my job, but I am sort ... ; Apprehensive	Skill context from ConvAI2 I have 3 children; I am a karate black belt; ... Skill context from WoW Pecan pie; Pecan pie is a pie of pecan nuts mixed ... Skill context from ED None
Dialogue A: Do you have much experience using the different types of cleats? Which do you like best? B: I have a little. I also know about air jordans, a brand of footwear also popular with athletes. A: I enjoy air jordans as well. I like to play soccer, and it's extremely easy to get good footwear with or without shoes. B: I agree. Air jordans are good for the environment and for calves, and it's really easy to pull off. A: Air jordans are generally made with the most material so it makes sense they'd be easy to use. B: And now shoes can go together with clothes as well, like any other type of material. A: Do you know much about aeros then? They have excellent fit and beauty. B: I have a fan and a hard time believing it but they are great shoes. A: I hope you can find a fan that is compatible with your shoes. B: I like my new balance flying fan. Do you have any?	

Table 14: Dialogue examples randomly sampled from BSBT dataset.

Example 5	
Skill context from ConvAI2 I drive an SUV; I've 3 dogs; I have a law degree; ...	Skill context from ConvAI2 I've 3 small dogs; I tutor kids in my spare time; ...
Skill context from WoW Mexico City	Skill context from WoW Mexico City; Mexico City is the capital and most ...
Skill context from ED New apartment, family and friends I love, this is ...; Content	Skill context from ED None
Dialogue	
A: Ok, that's nice. How modern is your capital city?	
B: It's been around a minute, capital and most populous city is Algiers, it's located in the far north of the country.	
A: Very cool. How do you feel about it? Do you like your city or your country?	
B: I love it. and love the people. What do you do for a living?	
A: I currently have a ba in law. Do you have any pets?	
B: I have 3 little ones, and 2 girls.	
A: Nice! Do you like them? Will you move?	
B: Yes , I would. The cost of living in Algiers is so high that I do not know if I would be able to afford it.	
A: It is expensive. Do you ever vacation? I like to visit the mountains.	
B: Yes. I visit a lot of places there.	
Example 6	
Skill context from ConvAI2 My favorite band is bon over; I am from Seattle; ...	Skill context from ConvAI2 I like to attend wine tours; I enjoy visiting museums; ...
Skill context from WoW Marduk (band)	Skill context from WoW Marduk (band); Marduk is a black metal band from ...
Skill context from ED My best friend's husband cheated on her with ...; Disgusted	Skill context from ED None
Dialogue	
A: I have never heard of a chupacabra. Would you like to live forever?	
B: No. However I would like to listen to Oasis' hit "Live Forever" since you mention the idea.	
A: That is true, I was hoping I could be of some help.	
B: Anything that helps would be great . thank you very much	
A: You're welcome. There is also a band called ctr that plays a bit of rap as well.	
B: I'll have to give them a listen.	
A: There is an album you should check out. It's called the on the court, you should really listen to it.	
B: I'll have to take a listen to it.	
A: I have to go now, but I hope you can tell them how much you love it and how you felt in this situation.	
B: Yes, I will do that as soon as possible.	

Table 15: Dialogue examples randomly sampled from BS \mathbb{B} T dataset.

Example 7	
Skill context from ConvAI2 I love coffee with cream; I went to new york city once; ...	Skill context from ConvAI2 I tutor kids in my spare time; I am very religious; ...
Skill context from WoW Dog	Skill context from WoW Dog; The domestic dog is a member of "Canis", ...
Skill context from ED We were living in Hawaii ...; Content	Skill context from ED None
Dialogue	
A: That is funny, what does jaja mean?	
B: Same as ll but in spanish!	
A: well, that is interesting, I guess. Although I don't use spanish much.	
B: Well it's used for many cultures, spanish comes after. What we get from it is also used in most western countries.	
A: Okay, that is cool. Thank you again, I am learning spanish also!	
B: I'm learning english as well! I'm so happy for you!	
A: Thanks, that is interesting. Its like learning a foreign language.	
B: That's a really cool skill to have. I bet you get to be learning all sorts of different languages.	
A: I sure do, because it's an important skill.	
B: Did you take it in high school? I know a lot of native spanish speakers.	
Example 8	
Skill context from ConvAI2 I like to make cookies; I like bagels with everything on; ...	Skill context from ConvAI2 I love the smell of beans; ...
Skill context from WoW Obesity in the United States	Skill context from WoW Obesity in the United States; Obesity in the United ...
Skill context from ED My girlfriend told me she's pregnant; Surprised	Skill context from ED None
Dialogue	
A: Oh wow what kinds of videos do you have?	
B: I do mainly make up tutorials. Do you have any side hobbies?	
A: I like making cookies and eating cookie dough, I love bagels too.	
B: I love bagels as well! What else do you like to do?	
A: Yum i love bagels too! And I also love playing sports, but I'm too busy at this moment.	
B: You should try joining a fitness team! It will make you feel better about you being busy.	
A: lol I'll, but I'd feel like i wouldn't be ready.	
B: You could always try it out! It's always fun to try out new things!	
A: I would if I could I really appreciate new things and learn new things from people like you.	
B: Sounds like you have an awesome hobby! Thanks for chatting.	

Table 16: Dialogue examples randomly sampled from BSBT dataset.