

LAVA: Latent Action Spaces via Variational Auto-encoding for Dialogue Policy Optimization

Nurul Lubis, Christian Geishaus, Michael Heck, Hsien-chin Lin,
Marco Moresi, Carel van Niekerk and Milica Gašić

Heinrich Heine University Düsseldorf, Germany

{lubis, geishaus, heckmi, linh, moresi, niekerk, gasic}@hhu.de

Abstract

Reinforcement learning (RL) can enable task-oriented dialogue systems to steer the conversation towards successful task completion. In an end-to-end setting, a response can be constructed in a word-level sequential decision making process with the entire system vocabulary as action space. Policies trained in such a fashion do not require expert-defined action spaces, but they have to deal with large action spaces and long trajectories, making RL impractical. Using the latent space of a variational model as action space alleviates this problem. However, current approaches use an uninformed prior for training and optimize the latent distribution solely on the context. It is therefore unclear whether the latent representation truly encodes the characteristics of different actions. In this paper, we explore three ways of leveraging an auxiliary task to shape the latent variable distribution: via pre-training, to obtain an informed prior, and via multitask learning. We choose response auto-encoding as the auxiliary task, as this captures the generative factors of dialogue responses while requiring low computational cost and neither additional data nor labels. Our approach yields a more action-characterized latent representations which support end-to-end dialogue policy optimization and achieves state-of-the-art success rates. These results warrant a more wide-spread use of RL in end-to-end dialogue models.

1 Introduction

With the rise of personal assistants, task-oriented dialogue systems have received a surge in popularity and acceptance. Task-oriented dialogue systems are characterized by a user goal which motivates the interaction, e.g. booking a hotel, searching for a restaurant, or calling a taxi. The dialogue agent is considered successful if it is able to fulfill the user goal by the end of the interaction. Traditionally, a dialogue system is built using the divide and conquer approach, resulting in multiple modules that together form a dialogue system pipeline: a natural language understanding (NLU) module, a dialogue state tracker (DST), a dialogue policy, and a natural language generation (NLG) module. Each module has well-defined input and output, and can be trained using machine learning (Young et al., 2010; Thomson and Young, 2010; Williams, 2006; Henderson et al., 2013) provided that adequately labeled data is available. However, there is a loss of information between the modules. The availability of powerful deep learning methods has recently led to a surge in end-to-end training approaches (Bordes and Weston, 2017; Wen et al., 2017; Madotto et al., 2018), which aim to map user utterances directly to responses in a sequence-to-sequence fashion.

Two fundamental properties of task-oriented systems are the ability to remember everything that is important from the conversation so far – *tracking*, and the ability to produce a response that steers the conversation towards successful task completion – *planning*. Within the modular approaches the role of tracking is taken by the DST, optimized using supervised learning (SL), while the role of planning is taken by the dialogue policy, optimized via reinforcement learning (RL). Unlike their modular counterparts, end-to-end systems typically only deploy SL, which relies on language modeling techniques that

directly optimize the likelihood of the data under the model parameters, neglecting planning altogether. This line of research has hugely benefited from large pre-trained transformer-based models such as BERT and GPT-2 (Devlin et al., 2019; Radford et al., 2019; Hosseini-Asl et al., 2020; Peng et al., 2020).

Only few works in end-to-end task-oriented dialogue systems deploy RL (Mehri et al., 2019b; Zhao et al., 2019). Word-level RL views each word of the entire system vocabulary as an action in a sequential decision making process. This blows up the action space size and the trajectory length, hindering effective learning and optimal convergence. The challenge of credit assignment and reward signal propagation is further compounded by the typically sparse rewards in dialogue. Last but not least, simultaneously optimizing language coherence and decision making within one model can lead to divergence. Thus, effectively incorporating RL in the end-to-end setting remains a challenge.

In the recently proposed latent action reinforcement learning (LaRL), a latent space between the context encoder and the response decoder serves as action space of the dialogue agent (Zhao et al., 2019). Decoding responses conditioned on the latent variable has the benefit of decoupling action selection and language generation, as well as shortening the dialogue trajectory, leading to improved performance. However, this approach optimizes the latent space using an uninformed prior without taking into consideration the actual distribution of the responses. Furthermore, the latent space model is conditioned only on the context. Therefore it is unclear whether the latent variables truly encode the characteristics of different dialogue actions, or whether it encodes the dialogue context instead. Because RL optimizes action selection and planning, it is important that it is performed on an action-characterized space.

In this paper, we propose an unsupervised approach for optimizing the latent action representation for end-to-end dialogue policy optimization with RL. Our contributions are as follows:

- We propose to optimize latent representations to be action-characterized. Action-characterized representations encode similar actions close to each other and allow interpolation of actions. This leads to a more practical and effective end-to-end RL.
- We explore three ways of leveraging an auxiliary task to shape the latent variable distribution; via pre-training, to obtain an informed prior, and via multitask learning. As auxiliary task, we choose response auto-encoding, as this captures generative factors of the dialogue responses. Unlike contemporary transformer-based approaches, this requires no additional data and has low computational cost. Our analysis shows that the learned latent representations encode action-relevant information.
- We show that our approach achieves state-of-the-art match and success rates on the multi-domain MultiWoZ 2.0¹ (Budzianowski et al., 2018).

This work acts as a proof of concept that we can induce an action-characterized latent space in an unsupervised manner to facilitate more practical and effective RL. The overall performance could likely be improved further by using more sophisticated encoding and decoding models, but this goes beyond the scope of this work. The results obtained here already warrant a more wide-spread use of RL in end-to-end dialogue models.

2 Related Work

Research in end-to-end task-oriented systems is largely inspired by the success of sequence-to-sequence modeling for chat-oriented systems (Serban et al., 2016). Representation learning has been shown to be useful for end-to-end systems, allowing more effective information extraction from the input. A common method leverages pre-training objectives that are inspired by natural language processing tasks, e.g. next-utterance retrieval (Lowe et al., 2016) and generation (Vinyals and Le, 2015). Naturally, this requires sufficient amounts of additional data, and often labels. The choice of the pre-training objective has been demonstrated to highly influence generalizability of the learned representation (Mehri et al., 2019a).

More recently, researchers have also investigated representation learning towards a better modeling of dialogue response. For example, Zhao et al. (2020) have investigated the use of language modeling tasks

¹The codebase is accessible at: <https://gitlab.cs.uni-duesseldorf.de/general/dsml/lava-public>. We also achieve state-of-the-art results on MultiWoZ 2.1

such as masking and sequence ordering for response generation. With variational models, latent space that spans across domains can be induced using dialog context-response pairs as well as a set of response-dialog act pairs (Zhao and Eskenazi, 2018). Similarly, dialogue context preceding and succeeding a response can be used in a skip-thought fashion to train response embeddings (Zhao et al., 2018). It has been reported that such representations allows few-shot domain adaptation using only raw dialog data (Shalymov et al., 2019). State labels and their transitions have also been utilized for learning action representations (Huang et al., 2020). Performing RL on top of the learned latent variable space has been shown to lead to a better policy compared to word-level RL, due to the condensed representation and shorter trajectory (Zhao et al., 2019). While improvement on metrics such as task success and entity recognition are reported, lack of interpretability and controllability remains a major challenge in this family of models.

3 Preliminaries

In an end-to-end framework, dialogue policy optimization with RL typically consists of two steps: SL and policy gradient RL. In the SL step, the model learns to generate a response x based on some dialogue context c , updating its parameters θ to maximize the log likelihood of the data,

$$L_{\text{SL}}(\theta) = \mathbb{E}_{x,c}[\log p_{\theta}(x|c)]. \quad (1)$$

Subsequently, the RL step updates the model parameter w.r.t. the task-specific goal, reflected as a reward. In a dialogue with T steps, for a specific time-step t , immediate reward r_t , and discount factor $\gamma \in [0, 1]$, the discounted return is defined as $R_t = \sum_{k=t}^T \gamma^{k-t} r_k$. The model tries to maximize the expected return from the first time-step onwards, written as $J(\theta) = \mathbb{E}_{\theta}[\sum_{t=0}^T \gamma^t r_t]$.

In word-level RL, every output word is treated as an action step, yielding the following policy gradient:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\theta} \left[\sum_{t=0}^T \sum_{j=0}^{U_t} R_{tj} \nabla_{\theta} \log p_{\theta}(w_{tj} | w_{<tj}, c_t) \right] \quad (2)$$

where T is the total number of turns in the dialogue, U_t is the total number of tokens in the response at turn t and j is the index of each token w . R_{tj} denotes the discounted return of the j -th token at turn t . In this policy gradient form, the action space is the vocabulary size of the system $|V|$, and the trajectory length is $\sum_{t=0}^{T-1} U_t$, making RL in this space extremely challenging.

The introduction of a latent variable z allows us to factorize the conditional distribution into $p(x|c) = p(x|z)p(z|c)$. By treating the latent space z as the action space, the action space size and trajectory length are reduced (Zhao et al., 2019). For a dialogue with T turns, policy gradient is now given by

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\theta} \left[\sum_{t=0}^T R_t \nabla_{\theta} \log p_{\theta}(z_t | c_t) \right], \quad (3)$$

where R_t denotes the discounted return at turn t . Zhao et al. (2019) have examined two types of latent variables; categorical and continuous. The categorical latent variable takes form as M independent K -way categorical random variables, while the continuous one is modeled as M dimensional multivariate Gaussian distribution with a diagonal covariance matrix. The latent variable distribution is learned during the SL step using stochastic variational inference by maximizing the evidence lowerbound (ELBO) – the lowerbound on the data log likelihood,

$$L_{\text{full}}(\theta) = \mathbb{E}_{q_{\theta}(z|x,c)}[\log p_{\theta}(x|z)] - \text{D}_{\text{KL}}[q_{\theta}(z|x,c) || p_{\theta}(z|c)]. \quad (4)$$

To combat exposure bias, Zhao et al. (2019) introduced a “lite” version of ELBO by assuming the posterior $q_{\theta}(z|x,c)$ to be the same as the encoder $p_{\theta}(z|c)$. This eliminates the second term of the ELBO objective. To counter overfitting, the posterior is regularized with some weight β to be similar to certain priors, in this case a uniform distribution for categorical latent variables, or a normal distribution for continuous ones. The “lite” ELBO objective has been shown to outperform the full ELBO objective, and is written as:

$$L_{\text{lite}}(\theta) = \mathbb{E}_{p_{\theta}(z|c)}[\log p_{\theta}(x|z)] - \beta \text{D}_{\text{KL}}[p_{\theta}(z|c) || p(z)]. \quad (5)$$

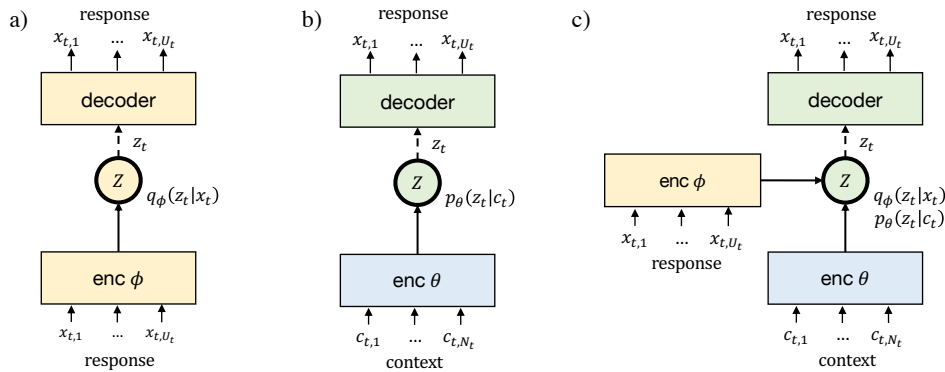


Figure 1: We aim to train an action-characterized distribution of latent variables Z to support RL in an end-to-end setting. a) VAE pre-training, b) new encoder is initialized and connected to the pre-trained Z and VAE decoder. Overall fine-tuning optimizes the entire network (LAVA_ptA), while selective fine-tuning optimizes only the new encoder (LAVA_ptS), c) new encoder is initialized and used in tandem with VAE encoder to obtain informative prior (LAVA_kl). The same architecture is also used for multitask learning (LAVA_mt), where we optimize for both tasks at the same time from scratch.

4 Latent Action Space via Auxiliary Task

As evidenced in Equation (3), to facilitate an effective learning it is important that the policy is trained using latent variables z that meaningfully represent actions. However, because existing methods model the distribution of z w.r.t. the context c (Equation (5)), we suspect that the latent variable is rather a representation of the state space instead of the action space. Furthermore, the distribution is regularized using an uninformed prior and without taking into account the actual distributions regarding dialogue responses for a given context. To induce a more action-like latent representations, we train the model on an auxiliary task that requires knowledge of responses to perform. We chose response auto-encoding (AE) as the auxiliary task using the variational auto-encoding (VAE) model. That is, given a response x we train the model to reconstruct the response via a latent space between the encoder and decoder (Figure 1a). With an uninformed prior $p(z)$, the pre-training objective for a set of parameters ϕ is:

$$L_{\text{vae}}(\phi) = \mathbb{E}_{q_{\phi}(z|x)}[\log p_{\phi}(x|z)] - D_{\text{KL}}[q_{\phi}(z|x)||p(z)]. \quad (6)$$

VAE models have been shown to be able to capture generative aspects of the samples they are trained on, resulting in good interpolation between latent variables (Kingma and Welling, 2014; Bowman et al., 2016). By training a VAE on dialogue responses, we aim to capture generative aspects of responses such as intent and domain information in an unsupervised manner.

We propose to utilize the VAE latent representations to condition dialogue systems to map encoded dialogue states to latent actions, instead of learning latent representations of the dialogue states. We call this approach **LAVA** (Latent Action via VAE). We explore three ways of leveraging the AE task to induce action-characterized latent variable distributions: as pre-training, as informed prior, and in a multitask learning fashion – pictured in Figure 1. Note that it is possible to swap the AE task with other tasks that target representation learning on the dialogue responses. In this work we utilize simple recurrent models as encoder and decoder to highlight the role of the latent dialogue action space. This allows us to pin any observed improvements on the latent space and dialogue policy. Other parts of the end-to-end dialogue system framework, such as encoding and decoding, are not within the scope of this work.

4.1 Auxiliary Task as Pre-training

The first method utilizes the auxiliary AE task to pre-train the latent representation and decoder (Figure 1b). Since the vocabularies of user and system turns vastly differ, a new dialogue system encoder for response generation (RG) is initialized and the VAE encoder is discarded. The new encoder is connected to the latent space and the decoder of the VAE. We experiment with two kinds of fine-tuning scheme:

overall (LAVA_ptA) and selective (LAVA_ptS). With LAVA_ptA, we update all parameters of the network during fine-tuning. On the other hand, LAVA_ptS blocks the gradient propagation to the latent space and the decoder and exclusively trains the encoder. This is equivalent to utilizing the generation modules from the VAE to serve as the NLG module, encouraging the model to focus on encoding the dialogue state (similar to a state tracker) and mapping it to a corresponding latent action (similar to a policy) during RG training. The loss function for both LAVA_ptA and LAVA_ptS is the “lite” ELBO objective (Equation 5).

4.2 Auxiliary Task as an Informative Prior

Secondly, we explicitly train the RG latent variable distribution to be close to that of the VAE. We call this set up LAVA_kl. As before, we start with a pre-trained VAE and a newly initialized RG encoder. We exploit the pre-trained latent representation in a novel way; the VAE encoder is not discarded and instead used in tandem to obtain an informed prior of the target response, pictured in Figure 1c. We use the latent distribution conditioned on the target $q_\phi(z|x)$ in the KL term penalty, replacing the uninformed prior used in previous works. This grounds the RG latent variable distribution to that of the VAE by penalizing divergence, while still optimizing it to fit the dialogue contexts. All parameters except the VAE encoder are further optimized with the following ELBO:

$$L_{\text{LAVA_kl}}(\theta) = \mathbb{E}_{p_\theta(z|c)}[\log p_\theta(x|z)] - \beta \text{D}_{\text{KL}}[p_\theta(z|c)||q_\phi(z|x)]. \quad (7)$$

4.3 Multitask Training between Main and Auxiliary Tasks

Lastly, we train a model to solve the RG and AE tasks in a multitask fashion, where RG is considered as the main task and AE as the auxiliary task. Multitask learning aims to improve learning efficacy by having a model solve multiple tasks at once, exploiting similarities across tasks (Caruana, 1997). Recent works have shown that dialogue system tasks also benefit from multitask learning (Rastogi et al., 2018; Zhu et al., 2019). RG and AE tasks are similar in that both aim to generate dialogue responses x , but they differ in the context they consider, also called the many-to-one multitask setting (Luong et al., 2015). RG attempts to generate the target response x given a dialogue context c , and AE tries to perform reconstruction given a response x . The two tasks share the latent space and decoder with a set of parameters ω but with separate encoders for RG and AE, with parameters θ and ϕ , respectively (Figure 1c). The aim is that the latent representation encodes more action-characterized features, since these are the common information required to fulfill both tasks. The two tasks are trained in an alternate fashion with an A:B ratio, i.e. for every A iterations of the main task, we train with the auxiliary task for B iterations. Unlike the previous methods, in multitask learning we start with a newly initialized model without pre-training. Each encoder receives an update only from its corresponding task, while the latent representation and decoder are trained on both tasks. The ELBO objectives are

$$L_{\text{LAVA_mt}}^{\text{RG}}(\omega, \theta) = \mathbb{E}_{p_\theta(z|c)}[\log p_\omega(x|z)] - \beta \text{D}_{\text{KL}}[p_\theta(z|c)||p(z)], \quad (8)$$

$$L_{\text{LAVA_mt}}^{\text{AE}}(\omega, \phi) = \mathbb{E}_{q_\phi(z|x)}[\log p_\omega(x|z)] - \beta \text{D}_{\text{KL}}[q_\phi(z|x)||p(z)]. \quad (9)$$

5 Experiment Setup

5.1 Corpus, Task, and Training Setup

We use the MultiWOZ 2.0 corpus (Budzianowski et al., 2018) to test the performance of the models. MultiWOZ is a collection of conversations between humans in a Wizard-of-Oz fashion, where one person plays the role of a dialogue system and the other one a user. The user is tasked to find entities, e.g. a restaurant or a hotel, that fit certain criteria by interacting with the dialogue system. The corpus simulates a multi-domain task-oriented dialogue system interaction, i.e. multiple domains may occur in the same dialogue or even the same turn. The corpus is fully annotated with a total of 10438 dialogues in English, it is one of the most challenging and largest corpora of its kind. We use the training, validation, and test set partitions provided in the corpus, amounting to 8438 dialogues for training, and 1000 each for validation and testing. All numbers reported are based on evaluation on the test set.

We aim to train a latent action representation that is effective for optimizing dialogue policies with RL in an end-to-end setting. This goal is best reflected by completion of the underlying dialogue task, measured in dialogue-level match and success rates. Match rate computes whether the criteria informed by the user (informable slots) are matched by the system, and success rate computes whether information requested by the user (requestable slots) are provided by the system. Match is a pre-requisite for a successful dialogue. For a long time the research in dialogue policy has only looked at success rates and user satisfaction (Lee and Eskénazi, 2012; Ultes et al., 2017), but as the line between policy and NLG becomes blurred we see the introduction of metrics such as BLEU and perplexity. However, these have been labeled early on to be potentially misleading as they correlate poorly with human judgement (Stent et al., 2005; Liu et al., 2016). Although we also report BLEU score for completeness, note that our methods are not targeted at improving BLEU.

We examine two tasks: 1) “dialogue context-to-response generation,” that is to generate the next dialogue response given dialogue history, as well as oracle dialogue state and database pointer from the corpus. The dialogue state is a binary vector representation of the user goal as inferred from the beginning of the dialogue up to the current turn. On the other hand, the database pointer is a binary vector representation of the entity matching the criteria in the dialogue state. 2) “End-to-end modeling” which takes only dialogue history for response generation. Works in end-to-end modeling typically utilize intermediate models in the pipeline to predict labels such as dialogue state and database pointer. In this work we utilize the latent action in an end-to-end fashion without the use of any intermediate labels, encouraging the model to fully exploit the latent variables. All experiments are conducted on delexicalized dialogues, where occurrences of slot values are replaced with their corresponding slot tokens, for example “in southern part of town” becomes “in [value_area] part of town.”

Our training consists of two steps: techniques presented in Section 4, followed with RL using REINFORCE. For a fair comparison with existing works, we adopt the novel RL setup proposed by Zhao et al. (2019): 1) For each RL episode, sample a dialogue from the corpus. 2) Run the model to generate a response for each system turn. However, the next user turn in the dialogue is not altered and simply retrieved from the corpus. 3) Compute success rate of the dialogue based on system response and use this as reward signal to compute policy gradient and update model parameters (Equation 3).

5.2 Model

Our primary focus are the latent representations induced by the proposed methods, their effect on reinforcement learning and the final performance of the model. To highlight the role of the latent dialogue action space, we use simple recurrent models as encoder and decoder for both the VAE and the dialogue system. We limit the model vocabulary to the most frequent 1000 tokens. We truncate the dialogue history to the last 2 turns for the context-to-response task, and the last 4 turns for the end-to-end modeling task. For both AE and RG tasks, the encoder is a GRU-LSTM (Cho et al., 2014) with attention and size 300 which outputs a vector with size 600. We tested both categorical and continuous latent variables. The categorical latent space is 10 independent 20-way categorical variables ($M = 10$, $K = 20$) and the continuous space is set to size $M = 200$. In the categorical case, the decoder is of size 150 with attention, and in the continuous case the decoder is of size 300. In choosing the hyperparameters, we follow the experimental set up reported by Zhao et al. (2019). We tested a few different set-ups, for example by varying the size of the latent space or the network, but found that the reported settings give the most optimal performance. One exception is the weight β for the KL term, which we set to 0.01 for all models other than LAVA_kl, which performs best with 0.1. Multitask training ratio is set to 10:1.

Unlike transformer-based architectures, the training of our models are computationally light and fast. One training takes between 1-3 hours using a single RTX 2080 GPU. While typical SL training requires 80-85 epochs with batch size of 128, LAVA_kl converges in under 20 epochs. LAVA_ptA, LAVA_ptS, and LAVA_mt convergence varies at around 20-50 epochs. LAVA models likely benefit from the pre-trained VAE models and is therefore able to converge faster.

Model	SL			+RL		
	match	succ.	BLEU	match	succ.	BLEU
LiteAttnCat*	65.77	57.26	0.18	83.68	78.18	0.12
LAVA_ptA_cat	70.57	58.56	0.18	85.09	77.08	0.12
LAVA_ptS_cat	64.56	54.65	0.19	83.48	79.87	0.12
LAVA_kl_cat	71.97	57.96	0.18	85.59	83.38	0.12
LAVA_mt_cat	60.46	51.05	0.18	84.88	80.98	0.10
LiteGauss*	68.27	57.06	0.19	75.88	66.47	0.14
LAVA_ptA_gauss	64.46	54.55	0.19	77.78	62.26	0.15
LAVA_ptS_gauss	67.37	53.85	0.18	77.78	61.76	0.13
LAVA_kl_gauss	68.87	58.66	0.19	79.28	64.96	0.09
LAVA_mt_gauss	59.42	49.33	0.19	82.78	70.37	0.14
Seq2Seq	58.66	52.25	0.20	81.58	75.18	0.15
SFN (Mehri et al., 2019b)	65.80	51.30	0.17	82.70	72.10	0.16
LiteAttnCat (Zhao et al., 2019)	67.97	57.36	0.19	82.80	79.20	0.12

Table 1: Best performance of our proposed methods in comparison with *reproduced and reported base-lines methods that employ RL. Our best model LAVA_kl_cat surpasses the baselines in both match and success rates.

Model	Match	Success	BLEU	Transformer	RL
Human	90.40	82.28	-	-	-
SimpleTOD (Hosseini-Asl et al., 2020)	88.90	67.10	0.16	✓	-
ARDM (Wu et al., 2019)	87.40	72.80	0.20	✓	-
DAMD (Zhang et al., 2020)	89.20	77.90	0.18	✓	-
SOLOIST (Peng et al., 2020)	89.60	79.30	0.18	✓	-
MarCo (Wang et al., 2020b)	92.30	78.60	0.20	✓	-
HDNO (Wang et al., 2020a)	96.40	84.70	0.18	-	✓
LAVA_kl_cat + RL (ours)	97.50	94.80	0.12	-	✓

Table 2: Comparison of our best performing model with existing works on the same task. For a fair comparison we adjust the performance of our best model by recalculating the match and success rates using a modified script released in the MultiWoZ repository². We also note whether the approaches employ a transformer-based architecture and RL.

6 Experimental Results

6.1 Context-to-Response Generation

Table 1 presents the performance of our models in comparison with 1) sequence-to-sequence (Seq2Seq) and structured fusion network (SFN) as baseline models that do not employ latent variables and 2) LaRL models, LiteAttnCat and LiteGauss, which we reproduced using the public code (marked with *) (Zhao et al., 2019). The Seq2Seq model shows best BLEU score compared to any of our models, however its match and success rate are consistently lower than our categorical models. This is not surprising, as Seq2Seq is optimized only to maximize the likelihood of the data. With the categorical latent variable, training a model on top of the pre-trained VAE, either in a selective manner or not, improves dialog-level performance. Using the VAE as informed prior gives us improvements on both match and success rates while maintaining the BLEU score, resulting in our best performing model. Although we find that continuous latent space is not as effective for RL, the proposed multitask training still surpasses the LiteGauss baseline and gives the best performance when Gaussian latent space is utilized.

We also compare our best performing model with existing works tackling the same task. Unlike the baseline models in Table 1, these works are evaluated with a new evaluation script recently published in the official MultiWoZ repository. The new evaluation script differs to the original one in the treatment of one specific case in the train domain, where a train matching user requirement is found but the user does not request the train ID. The original script underestimates the model performance as the train ID is checked regardless, while the new evaluation script do not check this further. For a fair comparison with relevant state-of-the-art models, we re-compute the match and success rates of the model using the new evaluation script. The numbers are reported in Table 2.

Model	labels			dialog-level		turn-level
	dialogue state	DB Search	action	match	succ.	BLEU
DAMD (Zhang et al., 2020)	gen	oracle	gen	76.30	60.40	0.18
SimpleTOD (Hosseini-Asl et al., 2020)	gen	-	gen	84.40	70.10	0.15
SOLOIST (Peng et al., 2020)	gen	gen	-	85.50	72.90	0.16
LAVA_kl_cat + RL	-	-	latent	91.80	81.80	0.12

Table 3: Comparison to state-of-the-art models on end-to-end generation task. Labels can come from the corpus (oracle), prediction using supervised models (gen), or in our case from latent space (latent). Our model is able to perform well without any additional labels by solely leveraging the latent variable.

We reach state-of-the-art inform and success rates, surpassing that of existing works which tackle the same task and even human performance on the test set at 82.28%. This result demonstrates the advantage of RL in the end-to-end setting. Optimizing for reward allows us to reach a higher success rate compared to solely focusing on accurately generating the target response. Our reward definition only takes into account the dialogue success rate, however because match is a prerequisite of dialogue success we are able to harmoniously optimize both match and success rates during RL. It is also interesting to note that our best model has higher match to success ratio compared to existing works, i.e. , we are able to achieve success on most dialogue where match occurs, while existing works fail more often in providing user with the information they require even when match already occurs. Moreover, unlike state-of-the-art models, our end-to-end setting uses simple encoder-decoder models without explicit dialogue state tracking, and therefore it is evident that the improvements are contributed by the latent action space. Combining our methods with more powerful models would be straightforward and we expect it to further improve the performance.

Notwithstanding, with regards to the very high performance, it is important to note the limitation of the current evaluation set up, most importantly that the dialogue trajectory is only estimated, since the user turn is obtained from data regardless of system response. To better gauge the performance in real dialogue with humans, user evaluation needs to be conducted in the future.

6.2 End-to-End Generation

We also tested our latent action in a fully end-to-end fashion, i.e. without using dialogue state labels and database pointers. Unlike existing works that train intermediate models for predicting labels such as dialogue state and action, our aim is to rely solely on the latent variable for forming a dialogue policy. We take our LAVA_kl_cat as pre-trained model, and further perform SL and RL exclusively with raw dialogue data in an end-to-end setting.

We present the performance of our best model in comparison with existing works in Table 3, along with the types of labels they utilize in the pipeline. Performance of these models are also computed with the new evaluation script as in Table 2. Consistent with previous task, our model is outperforming the other models in terms of success and inform rates. The result confirms that our model is able to optimize its dialogue policy by leveraging action-relevant information that is encoded in the latent variable, e.g. action type and domain, even in an extreme setting where no additional label is utilized in the pipeline.

7 Latent Space Analysis

7.1 Clustering Metrics and Projection

We investigate whether the latent variables are grouped together according to true action or domain labels. One method is to quantify the quality of clusters that are formed in the latent space w.r.t. action and domain labels. We use the Calinski-Harabasz index, which measures the ratio of the sum of between-clusters dispersion and of inter-cluster dispersion for all clusters (Caliński and Harabasz, 1974). A higher Calinski-Harabasz score relates to a model with better defined clusters².

Table 4 compares the scores of the LAVA_kl and LAVA_mt methods with their corresponding reproduced LaRL baselines (Zhao et al., 2019). We observe that for all models domain labels are better clustered than action, which is expected because 1) the amount of unique domains is much smaller than

²We note that extremely high scores could signal cluster overfitting. See Appendix B for more details.

Model	Categorical				Gaussian			
	SL		RL		SL		RL	
	Domain	Action	Domain	Action	Domain	Action	Domain	Action
LaRL*	93.19	23.30	121.15	17.50	47.86	13.04	71.49	13.25
LAVA_kl	104.92	25.28	158.00	41.75	106.51	20.00	128.51	19.13
LAVA_mt	25.37	6.64	247.92	16.85	66.54	16.91	80.09	18.50

Table 4: Clustering metrics scores \uparrow . Our LAVA_kl_cat and LAVA_mt_gauss show harmonious improvement of domain and action clusters, as well as a nice balance between the domain and action scores.

that of actions, and 2) domain information is explicitly expressed at word level, while action is concerned with the intent of the utterance. Note that since domain information is part of the action label, efficient domain clustering is also beneficial for inferring actions. Performing RL on top of SL consistently improves the clustering scores. However, in some cases improvement on domain clusters comes at the cost of action cluster, e.g. categorical LaRL and LAVA_kl_gauss. We observe that good dialogue performance aligns with harmonious improvement of domain and action clusters, as well as a nice balance between the domain and action scores, as demonstrated by LAVA_kl_cat and LAVA_mt_gauss.

We visually assess the latent space by projecting the latent action of each input in the training set with t-SNE (Maaten and Hinton, 2008) and analyzing the cluster that formed w.r.t. domain and action labels. We compare our best proposed model LAVA_kl_cat with the baseline LiteAttnCat, before and after RL, presented in Figure 2. While LiteAttnCat sees good domain cluster definition after SL and RL, it loses some action cluster definition after RL. On the other hand, when training on the informed prior, we obtain clusters that are tighter and farther apart from each other. Performing RL on top of SL moves the clusters inwards and improves cluster definition without causing significant transformation of the latent space. This indicates that the proposed method allows the model to put more focus on learning a dialogue policy without having to radically modify the latent action representation. This also signals that the model is equipped with an effective action space since the beginning of RL, which boosts learning.

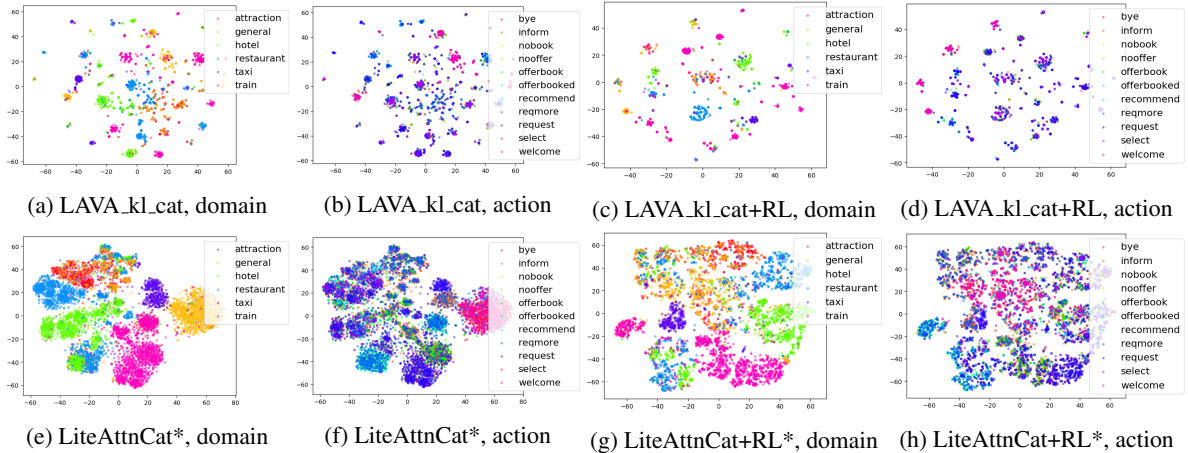


Figure 2: LiteAttnCat* and LAVA_kl_cat latent space projection before and after RL. While LiteAttnCat loses action cluster definition, our LAVA_kl_cat improves both domain and action clusters definition without causing significant latent space transformation. Higher resolution version is in Appendix C.

7.2 Latent Variable Traversal

Latent variable traversal can be employed to qualitatively analyze the relationship between different parts of the latent space. This is done by taking two points from the latent space, traversing the space between them by interpolating intermediate variables, and generating a sample for each variable. A latent space that meaningfully encode the generative factors would be able to produce reasonable samples with gradating similarities to the opposite ends of the traversal.

We select two dialogue contexts where the target responses perform similar actions but in different domains. Traversal of latent variables from LAVA_kl_cat and reproduced LiteAttnCat models for these

LiteAttnCat*
<p>yes, [hotel_name] is a [value_pricerange] -ly priced [value_count] star guesthouse located in the [value_area]. it is [value_pricerange] -ly priced and has [value_count] stars . would you like to book a room?</p> <p>sure, i can help you with that . what would you like to know?</p> <p>i have the [hotel_name] located at [hotel_address]. would you like me to book it for you?</p> <p>i have the [restaurant_name] located at [restaurant_address]. would you like to book a table?</p> <p>the address is [restaurant_address] and the phone number is [restaurant_phone].</p> <p>the [hotel_name] is a [value_count] star guesthouse in the [value_area] area . it is [value_pricerange] -ly priced and has free wifi and parking . would you like me to book it for you?</p> <p>i have booked you a taxi . it will be a [taxi_type] and the contact number is [taxi_phone].</p>
LAVA_kl.cat
<p>[hotel_name] is a guesthouse in the [value_area] area . it is [value_pricerange] -ly priced and has [value_count] stars . would you like to book a room?</p> <p>i would recommend [hotel_name]. it s a [value_count] star guesthouse in the [value_area]. would you like to book a room?</p> <p>the reference number for the train is [train_reference].</p> <p>i am sorry, i am not able to book that . i can book you a room at the [restaurant_name] if you would like.</p> <p>i am sorry, but i am unable to book it right now . is there anything else i can help you with?</p> <p>i can book that for you now.</p> <p>i have booked you a [taxi_type]. the contact number is [taxi_phone]. can i help you with anything else?</p>

Table 5: Latent variable traversal between two responses with dialogue actions related to booking in different domains. Traversal on LAVA_kl_cat shows smooth transition with consistent underlying action.

contexts is presented in Table 5. We observe that the proposed model is able to transition smoothly from an inform and offerbook action in the hotel domain, to providing booking confirmation in the taxi domain. On the other hand, latent representations trained without the support of VAE generate other actions in the traversal, signaling that the action features are not encoded effectively in the latent representation. It is also worthwhile to note that the two actions are closer to each other in the LAVA_kl.cat model, and farther apart in the LiteAttnCat model. With LAVA_kl_cat, traversal after RL yields the same responses while LiteAttnCat shows improved traversal. This echoes our previous analysis that the proposed method is equipped with an action-characterized action space since the beginning of RL, which supports effective and practical RL with end-to-end dialogue models.

8 Conclusion and Future Work

This work acts as proof of concept that we can induce action-characterized latent representations in an unsupervised manner to facilitate a more practical and effective RL with end-to-end dialogue models. We explore ways to obtain action-characterized latent representations via response variational auto-encoding, which captures generative aspects of responses. Treating these latent representations as actions allows effective optimization with RL. Unlike contemporary transformer-based approaches, our method requires no additional data and has low computational cost. We are able to achieve state-of-the-art success rate on the challenging MultiWoZ 2.0 corpus on both context-to-response generation as well as the end-to-end modeling task. Our analyses show that the proposed methods result in latent representations that cluster well w.r.t. domain and action labels, and encode similar actions close to each other. In this paper, we utilize simple recurrent models to highlight the merit of the proposed training methods, which means each component can be replaced with stronger models to further improve performance. We believe our method has high potential for end-to-end domain adaptation and offline policy learning with RL. We look forward to improve our work by utilizing longer context and performing RL in a stricter setting where the dialogue trajectory is more accurately estimated. We would also like to conduct human evaluation and analyze how our model performs in real dialogue interaction.

Acknowledgements

N. Lubis, M. Heck, and C. van Niekerk are supported by funding provided by the Alexander von Humboldt Foundation in the framework of the Sofja Kovalevskaja Award endowed by the Federal Ministry of Education and Research, while C. Geishauer, H-C. Lin and M. Moresi are supported by funds from the European Research Council (ERC) provided under the Horizon 2020 research and innovation programme (Grant agreement No. STG2018_804636). Computing resources were provided by Google Cloud.

References

- Antoine Bordes and Jason Weston. 2017. Learning end-to-end goal-oriented dialog. In *5th International Conference on Learning Representations*.
- Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. 2016. Generating sentences from a continuous space. *CoNLL 2016*, page 10.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Ultes Stefan, Ramadan Osman, and Milica Gašić. 2018. MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Tadeusz Caliński and Jerzy Harabasz. 1974. A dendrite method for cluster analysis. *Communications in Statistics-theory and Methods*, 3(1):1–27.
- Rich Caruana. 1997. Multitask learning. *Machine learning*, 28(1):41–75.
- Kyunghyun Cho, B van Merriënboer, Caglar Gulcehre, F Bougares, H Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Mihail Eric, Rahul Goel, Shachi Paul, Adarsh Kumar, Abhishek Sethi, Peter Ku, Anuj Kumar Goyal, Sanchit Agarwal, Shuyang Gao, and Dilek Hakkani-Tur. 2019. Multiwoz 2.1: A consolidated multi-domain dialogue dataset with state corrections and state tracking baselines. *arXiv preprint arXiv:1907.01669*.
- M. Henderson, B. Thomson, and S.J. Young. 2013. Deep neural network approach for the dialog state tracking challenge. In *Proceedings of SIGDIAL*.
- 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*.
- Xinting Huang, Jianzhong Qi, Yu Sun, and Rui Zhang. 2020. MALA: Cross-domain dialogue generation with action learning. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*.
- Diederik P Kingma and Max Welling. 2014. Auto-encoding variational bayes. *stat*, 1050:1.
- Sungjin Lee and Maxine Eskenazi. 2012. POMDP-based let’s go system for spoken dialog challenge. In *Proceedings of SLT*.
- Chia-Wei Liu, Ryan Lowe, Iulian Vlad Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2122–2132.
- Ryan Lowe, Iulian Vlad Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. On the evaluation of dialogue systems with next utterance classification. In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 264–269.
- Minh-Thang Luong, Quoc V Le, Ilya Sutskever, Oriol Vinyals, and Lukasz Kaiser. 2015. Multi-task sequence to sequence learning. *arXiv preprint arXiv:1511.06114*.
- Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9(Nov):2579–2605.
- Andrea Madotto, Chien-Sheng Wu, and Pascale Fung. 2018. Mem2seq: Effectively incorporating knowledge bases into end-to-end task-oriented dialog systems. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1468–1478.
- Shikib Mehri, Evgeniia Razumovskaia, Tiancheng Zhao, and Maxine Eskenazi. 2019a. Pretraining methods for dialog context representation learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3836–3845, Florence, Italy, July. Association for Computational Linguistics.

- Shikib Mehri, Tejas Srinivasan, and Maxine Eskenazi. 2019b. Structured fusion networks for dialog. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 165–177.
- Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayandeh, Lars Liden, and Jianfeng Gao. 2020. SOLOIST: Few-shot task-oriented dialog with a single pre-trained auto-regressive model. *arXiv preprint arXiv:2005.05298*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- Abhinav Rastogi, Raghav Gupta, and Dilek Hakkani-Tur. 2018. Multi-task learning for joint language understanding and dialogue state tracking. In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*, pages 376–384.
- Iulian V Serban, Alessandro Sordani, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2016. Building end-to-end dialogue systems using generative hierarchical neural network models. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- Igor Shalyminov, Sungjin Lee, Arash Eshghi, and Oliver Lemon. 2019. Few-shot dialogue generation without annotated data: A transfer learning approach. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 32–39.
- Amanda Stent, Matthew Marge, and Mohit Singhai. 2005. Evaluating evaluation methods for generation in the presence of variation. In *international conference on intelligent text processing and computational linguistics*, pages 341–351. Springer.
- B. Thomson and S. Young. 2010. Bayesian update of dialogue state: A POMDP framework for spoken dialogue systems. *Computer Speech and Language*, 24(4):562–588.
- Stefan Ultes, Pawel Budzianowski, Inigo Casanueva, Nikola Mrksic, Lina Maria Rojas-Barahona, Pei-Hao Su, Tsung-Hsien Wen, Milica Gasic, and Steve J Young. 2017. Domain-independent user satisfaction reward estimation for dialogue policy learning. In *INTERSPEECH*, pages 1721–1725.
- Oriol Vinyals and Quoc Le. 2015. A neural conversational model. *arXiv preprint arXiv:1506.05869*.
- Jianhong Wang, Yuan Zhang, Tae-Kyun Kim, and Yunjie Gu. 2020a. Modelling hierarchical structure between dialogue policy and natural language generator with option framework for task-oriented dialogue system. *arXiv preprint arXiv:2006.06814*.
- Kai Wang, Junfeng Tian, Rui Wang, Xiaojun Quan, and Jianxing Yu. 2020b. Multi-domain dialogue acts and response co-generation. *arXiv preprint arXiv:2004.12363*.
- Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, Milica Gašić, Lina M. Rojas-Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. 2017. A network-based end-to-end trainable task-oriented dialogue system. In *Proceedings of EACL*.
- J.D. Williams. 2006. *Partially Observable Markov Decision Processes for Spoken Dialogue Management*. Ph.D. thesis, University of Cambridge.
- Qingyang Wu, Yichi Zhang, Yu Li, and Zhou Yu. 2019. Alternating recurrent dialog model with large-scale pre-trained language models. *arXiv preprint arXiv:1910.03756*.
- Steve Young, Milica Gašić, Simon Keizer, François Mairesse, Jost Schatzmann, Blaise Thomson, and Kai Yu. 2010. The Hidden Information State model: A practical framework for POMDP-based spoken dialogue management. *Computer Speech and Language*, 24(2):150–174.
- Yichi Zhang, Zhijian Ou, and Zhou Yu. 2020. Task-oriented dialog systems that consider multiple appropriate responses under the same context. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*.
- Tiancheng Zhao and Maxine Eskenazi. 2018. Zero-shot dialog generation with cross-domain latent actions. In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*, pages 1–10.
- Tiancheng Zhao, Kyusong Lee, and Maxine Eskenazi. 2018. Unsupervised discrete sentence representation learning for interpretable neural dialog generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1098–1107.
- Tiancheng Zhao, Kaige Xie, and Maxine Eskenazi. 2019. Rethinking action spaces for reinforcement learning in end-to-end dialog agents with latent variable models. 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 1208–1218.

Yufan Zhao, Can Xu, and Wei Wu. 2020. Learning a simple and effective model for multi-turn response generation with auxiliary tasks. *arXiv preprint arXiv:2004.01972*.

Chenguang Zhu, Michael Zeng, and Xuedong Huang. 2019. Multi-task learning for natural language generation in task-oriented dialogue. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1261–1266.

Appendix

A Additional Results on MultiWoZ 2.1

For a more complete comparison with existing works, we report additional results of our best model. We trained and tested LAVA_kl_cat + RL with MultiWoZ 2.1 dataset (Eric et al., 2019), and computed the match and success rates of the model using a new evaluation script published in the official MultiWoZ repository³. MultiWoZ 2.1 includes corrections to the dialogue state labels and canonicalization of the slot values in the utterance, reducing the labeling error and typos introduced by human worker. As previously explained, the new evaluation script differs to the old one in the treatment of one specific case in the train domain, where a train matching user requirement is found but the user does not request the train ID. The old script underestimates the model performance as the train ID is checked regardless, while the new evaluation script do not check this further. The results are presented in Table 6. Training and testing with MultiWoZ 2.1 yields lower match and success rates, however BLEU score is improved. At the time of publication, our best model achieves state-of-the-art results.

Model	Match	Success	BLEU	Transformer	RL
SimpleTOD (Hosseini-Asl et al., 2020)	85.10	73.50	0.16	✓	-
MarCo (Wang et al., 2020b)	92.50	77.80	0.19	✓	-
HDNO (Wang et al., 2020a)	92.80	83.00	0.18	-	✓
LAVA_kl_cat + RL (ours)	96.39	83.57	0.14	-	✓

Table 6: Comparison of our best performing model with existing works on MultiWoZ 2.1 data. We also note whether the approaches employ a transformer-based architecture and RL.

B Cluster Overfitting and Its Relation to Calinski-Harabasz Index

For a dataset E of size n_E and a center c_E forming k clusters, with each cluster q of size n_q consisting of a centroid c_q and a set of points C_q , the Calinski-Harabasz index CH is computed as

$$CH = \frac{tr(B_k)}{tr(W_k)} \times \frac{n_E - k}{k - 1}, \quad (10)$$

where

$$W_k = \sum_{q=1}^k \sum_{x \in C_q} (x - c_q)(x - c_q)^T, \quad (11)$$

$$B_k = \sum_{q=1}^k n_q (c_q - c_E)(c_q - c_E)^T, \quad (12)$$

and $tr(B_k)$ and $tr(W_k)$ are the traces of the between-cluster dispersion and the inter-cluster dispersion, respectively. It is evident that the smaller the inter-cluster dispersion, the higher the score will be. However in our experiments we find that high score could be a sign of cluster overfitting in the latent space, that is when the clusters do not preserve enough variability of the target responses. For example, when each domain is grouped into one small cluster and each cluster is placed far from each other, scoring on domain labels may translate to extremely high values, yet this is not desired in practice as we would

³<https://github.com/budzianowski/multiwoz>

like the latent representation to capture within-cluster varieties as well, or maybe split one domain into several clusters to better distinguish other action-relevant information. Furthermore, overfitting w.r.t. domain label set typically means poor fit w.r.t. action labels. Figure 3 presents an example.

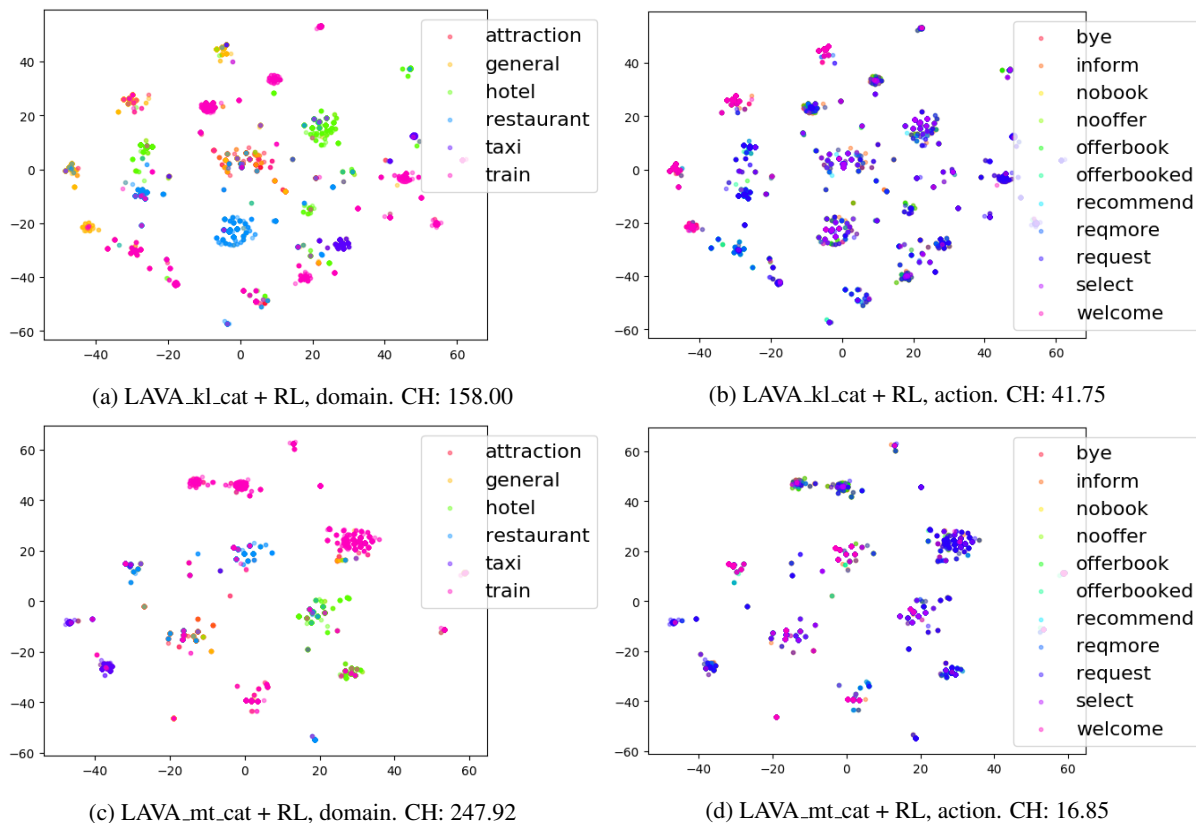


Figure 3: LAVA_kl_cat and LAVA_mt_cat latent space projection after RL and their respective Calinski-Harabasz scores (CH). All plots contain the same number of data points. Plots (c) and (d) show that LAVA_cat_mt + RL groups the datapoints into fewer clusters, and while this yields high domain CH score, the score for action clusters is low. On the other hand, although the domain CH is lower, LAVA_kl_cat shows better clustering fit for both domain and action.

C High Resolution Plots

We reproduce the cluster projections from Figure 2 in high resolution, presented in Figure 4.

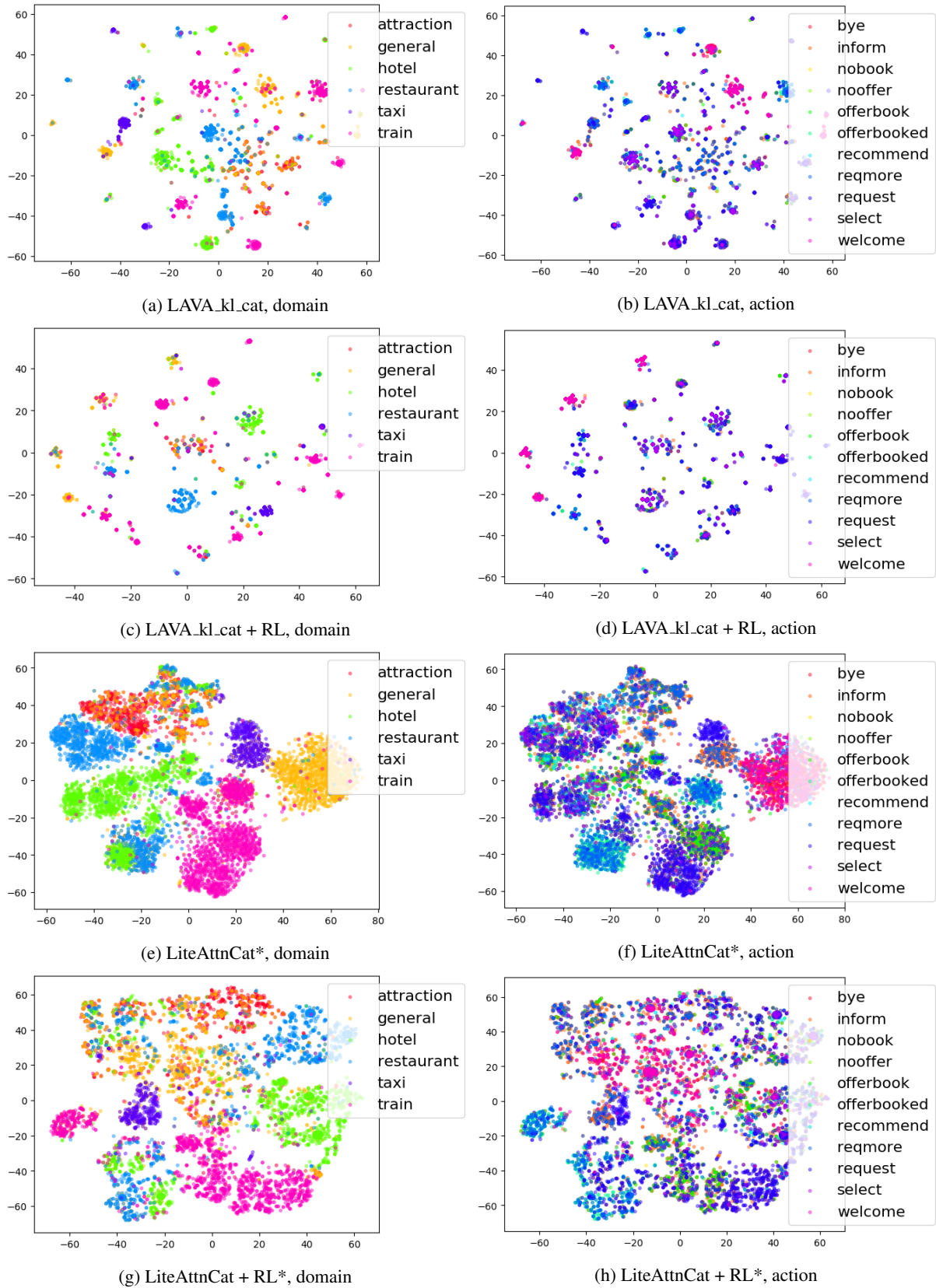


Figure 4: LiteAttnCat* and LAVA_kl_cat latent space projection before and after RL. While LiteAttnCat loses action cluster definition, our LAVA_kl_cat improves both domain and action clusters definition without causing significant transformation of the latent space.