

A Modular Task-oriented Dialogue System Using a Neural Mixture-of-Experts

Jiahuan Pei, Pengjie Ren, Maarten de Rijke
University of Amsterdam

Overview

- **Goal.** Generate a system response token by token given a dialogue context for Task-oriented Dialogue Systems (TDSs).
- **Problem.** Previous studies on end-to-end TDSs use a single-module model to generate responses for complex dialogue contexts. However, no model consistently outperforms the others in all situations.
- **Solution.**
 - ▷ Propose a neural Modular Task-oriented Dialogue System (MTDS) framework consisting of a chair and several experts.
 - ▷ Implement a Token-level Mixture-of-Expert (TokenMoE) model, where the experts predict multiple tokens at each timestamp and the chair determines the final generated token by fully considering the outputs of all experts.

System Response Generation in TDSs

- Given
 - ▷ a sequence of dialogue context $X = (x_1, \dots, x_m)$ with m tokens
 - ▷ and a sequence of system response $Y = (y_1, \dots, y_n)$ with n tokens
- The model aims to optimize the generation probability of Y conditioned on X , i.e., $p(Y|X)$.

MTDS Framework

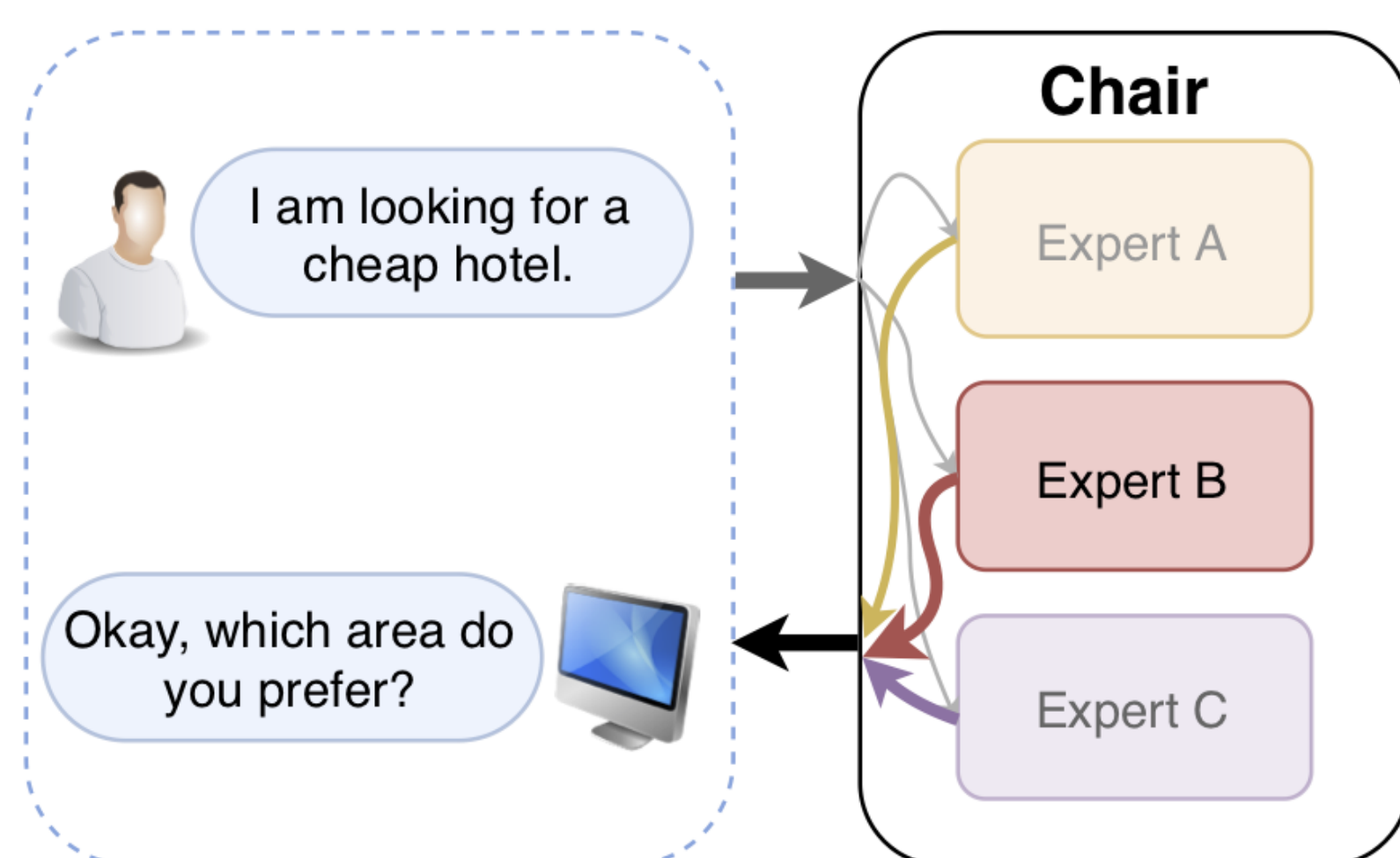


Figure: Modular Task-oriented Dialogue System (MTDS) framework.

- k **expert** bots, each of which is specialized for a particular *intent*.
- a **chair** bot, which learns to coordinate a group of expert bots to make an optimal decision.

Method

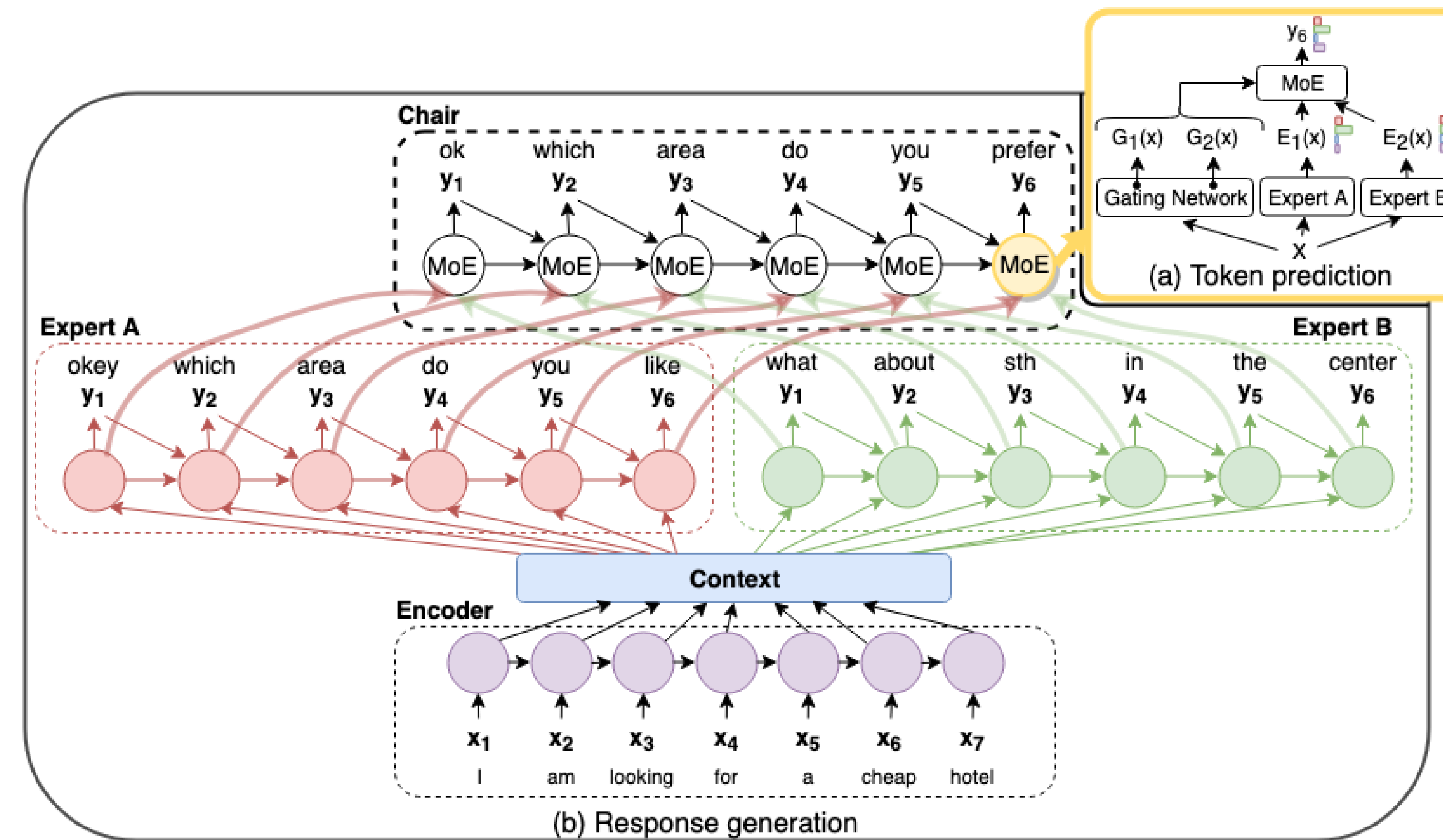


Figure: Overview of TokenMoE. Figure (a) illustrates how does the model generate the token y_6 given sequence X as an input. Figure (b) shows how does the model generate the whole sequence Y as a dialogue response.

Experimental setup

Research questions

- RQ1.** Is there a single model that consistently outperforms the others on all domains?
- RQ2.** Does the TokenMoE model outperform the state-of-the-art end-to-end single-module TDS model?
- RQ3.** How do the proposed token-level Mixture-of-Expert (MOE) scheme and the global-and-local learning scheme in the TokenMoE model affect the final performance?

We conducts experiments are on the MultiWOZ dataset and use three commonly used evaluation metrics:

- **Inform.** The fraction of responses that provide a correct entity out of all responses.
- **Success.** The fraction of responses that answer all the requested attributes out of all responses.
- **BLEU.** This is a score for comparing a generated response to one or more reference responses.

Results

Table: Performance of the single-module models vary on different domains.

	Inform (%)				Success (%)				BLEU (%)				Score			
	BSL	/V1	/V2	/V3	BSL	/V1	/V2	V3	BSL	/V1	/V2	/V3	BSL	/V1	/V2	/V3
Attraction	87.20	86.20	91.80	88.70	81.30	74.80	83.70	83.70	15.14	14.95	16.08	14.86	99.39	95.45	103.83	101.06
Hotel	89.90	93.90	89.90	90.30	87.50	91.70	87.40	89.10	16.60	15.60	15.11	14.13	105.30	108.40	103.76	103.83
Restaurant	89.20	91.70	86.40	86.10	85.80	87.80	84.00	83.40	17.07	17.70	16.07	17.34	104.57	107.45	101.27	102.09
Taxi	100.00	100.00	100.00	100.00	99.90	99.80	99.90	99.80	17.33	19.18	20.13	18.32	117.28	119.08	120.08	118.22
Train	77.70	77.70	79.00	81.60	75.60	74.80	77.20	79.60	20.35	15.64	22.81	20.62	97.00	91.89	100.91	101.22
Booking	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	22.05	21.61	21.96	22.06	122.05	121.61	121.96	122.06
General	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	20.21	19.53	20.13	20.80	120.21	119.53	120.13	120.80
UNK	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	12.40	11.75	13.12	11.80	112.40	111.75	113.12	111.80

Table: Comparison between TokenMoE, the benchmark baseline S2SAttnLSTM, and their variant models.

	Inform (%)	Success (%)	BLEU (%)	Score
S2SAttnLSTM	67.20	57.20	17.83	80.03
S2SAttnLSTM/V1	63.60	52.20	18.10	76.00
S2SAttnLSTM/V2	67.20	58.90	20.85	83.90
S2SAttnLSTM/V3	68.60	59.30	19.41	83.36
TokenMoE/V1	64.00	52.50	18.95	77.20
TokenMoE/V2	62.60	54.30	18.90	77.35
TokenMoE/V3	62.90	54.00	18.34	76.79
TokenMoE	75.30	59.70	16.81	84.31

Table: Comparison of TokenMoE with different learning schemes (S1, S2, S3, S4) and the benchmark baseline S2SAttnLSTM.

	Inform (%)	Success (%)	BLEU	Score
S2SAttnLSTM/V2	67.20	58.90	20.85	83.90
TokenMoE/S1	66.20	54.90	19.11	79.66
TokenMoE/S2	66.50	56.90	19.48	81.18
TokenMoE/S3	70.60	60.60	18.67	84.27
TokenMoE/S4	75.30	59.70	16.81	84.31

Table: An example of the generated responses of S2SAttnLSTM and TokenMoE. A user would prefer to get detail information of the train before booking a ticket.

Model	Response
S2SAttnLSTM	i have [value_count] trains that match your criteria . would you like me to book it for you ?
TokenMoE	i have train [train_id] that leaves at [value_time] and arrives at [value_time] . would you like me to book it ?

Conclusion

- A neural MTDS framework composed of a chair and several experts.
- A TokenMoE model under MTDS framework:
 - ▷ The experts make multiple token-level predictions at each timestamp.
 - ▷ The chair predicts the final generated token considering the whole outputs of all experts.
- Main findings:
 - ▷ No single-module TDS model can constantly outperform the others on all metrics for all cases.
 - ▷ TokenMoE outperforms the best single-module model by 8.1% of *Inform* rate and 0.8% of *Success* rate.
 - ▷ Learning scheme is an important factor of TokenMoE.



UNIVERSITY OF AMSTERDAM