

DYME: A Dynamic Metric for Dialog Modeling Learned from Human Conversations

Speaker: Florian von Unold

Authors: Florian von Unold^{*, 1, 2}, Monika Wintergerst^{*, 1}, Lenz Belzner³, Georg Groh¹

* These authors contributed equally

¹ Department of Informatics, Technical University of Munich, Boltzmannstr. 3, 85748 Garching, Germany

² MaibornWolff GmbH, Theresienhoehe 13, 80338 Munich, Germany

³ Faculty of Electrical Engineering and Information Technology, Technische Hochschule Ingolstadt,
Esplanade 10, 85049 Ingolstadt, Germany

> Code and models on GitHub: <https://github.com/florianvonunold/DYME>



Motivation

Context

Text-based chatbots

- Chatbots are a natural way of human-machine interaction
- Practical applications
 - Business: customer service
 - Healthcare: nursing, virtual dietary advice
 - Education: virtual teaching

→ Interested in improving open-domain chatbots

Motivation

Goal

Control data driven dialog generation to reflect human communication characteristics

Assumption: Controlling chatbots' „behavior“ allows for more effective communication (human-machine interaction)

Control **maybe** particularly important for dialog generation: output text highly sensitive to given dialog situation

Related Work

Fine-tune dialog models with Reinforcement Learning (Saleh et al., 2019; Xu, Wu, and Wu, 2018; Li et al., 2016)

- *Self-play / Self-talk*: two dialog models talk to each other to simulate conversations
- Optimize the model(s) w.r.t. certain reward functions (e.g., to get less repetitive models)
- Long-term influence of a generated utterance on the entire conversation is modeled
- Possibility to explore and move away from learned statistical dependencies

Related Work

Reward functions

- Learn reward functions from human preferences
- Developer-defined reward functions
- Learn reward functions from dialog data

Related Work

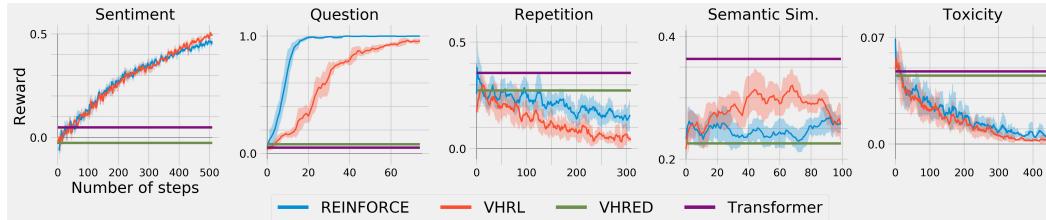
Reward functions

- Learn reward functions from human preferences
- Developer-defined reward functions
- Learn reward functions from dialog data

Related Work

Developer-defined reward functions

Variational Hierarchical Reinforcement Learning (VHRL) (Saleh et al., 2019)



4

Reward learning curves for the different metrics.

Transformer	Sentiment	Question	Repetition	Semantic Sim.	Toxicity
[Usr]: hello! [Bot]: i don't know. [Usr]: how are you? [Bot]: i don't know what that means.	[Usr]: hello! [Bot]: i'm in my sunroom a lot and it's a great feeling! [Usr]: sounds fun [Bot]: i'm excited!	[Usr]: hi there! [Bot]: what do you do for fun? [Usr]: i read books [Bot]: what do you like to do in your life?	[Usr]: how are you? [Bot]: i'm a reservoir engineer [Usr]: cool. i am still in college [Bot]: today is my first day of vacation	[Usr]: i like to listen to ed sheeran [Bot]: and what genre is it? [Usr]: pop [Bot]: any recommendations?	[Usr]: you're so stupid [Bot]: you're cute [Usr]: i hate you [Bot]: but i don't know why

4

Chat examples of the trained models after VHRL fine-tuning.

⁴ Figure and Table taken from (Saleh et al., 2019)

Related Work

Developer-defined reward functions

Variational Hierarchical Reinforcement Learning (VHRL) (Saleh et al., 2019)

$$reward_{VHRL} = 0.15 * sentiment + 0.25 * question + 0.5 * repetition + 0.05 * similarity + 0.05 * toxicity$$

According to the reward, every utterance at any given position in a dialog should adhere to these static, developer-defined metrics

Related Work

Developer-defined reward functions

Variational Hierarchical Reinforcement Learning (VHRL) (Saleh et al., 2019)

$$reward_{VHRL} = 0.15 * sentiment + 0.25 * question + 0.5 * repetition + 0.05 * similarity + 0.05 * toxicity$$

According to the reward, every utterance at any given position in a dialog should adhere to these static, developer-defined metrics
→ Not necessarily valid

Example dialog

[Speaker 1]: Hey!

[Speaker 2]: Hi, how are you?

[Speaker 1]: Fine, and you?

[Speaker 2]: I think about how sad it is that I've not yet been able to meet you in person... ← suitable utterance will receive lower reward

→ Possible problem: static rewards

Contributions

1. Analysis of metric dynamics in written human conversations → Metrics in human dialogs are dynamic.
2. **DYME**: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations → A dynamic metric can be learned from dialog data.

Experiments and Results

1. Analysis of metric dynamics in written human conversations

2. **DYME**: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Experiments and Results

Analysis of metric dynamics in written human conversations

Metric: any function that computes a real-valued number for a given utterance in a two-speaker, multi-turn conversation

Used metrics

- Adapted Reward Function Metrics⁵
 - Question, repetition, ...
- New Empathy-Based Metrics⁶
 - Emotional reaction, interpretation & exploration

Used datasets: human-written, multi-turn, two-speaker, open-domain dialogs

- DailyDialog (Li et al., 2017)
- EmpatheticDialogues (Rashkin et al., 2019)

⁵ Adapted versions of the reward functions in (Saleh et al., 2019)

⁶ Following the EPITOME framework of expressed empathy (Sharma et al., 2020)

Experiments and Results

Analysis of metric dynamics in written human conversations

- u_1 did you ever buy food from the snack stands near our hotel ?
- u_2 yes , several times .
- u_3 how do you like them ?
- u_4 not bad .
- u_5 i always have the temptation to eat something there .
- u_6 then , why did n't you do that ?
- u_7 i do n't know how much we can trust them . do you have any ideas ?
- u_8 some of them , i think , are not good .
- u_9 it does n't taste good ?
- u_{10} no , i mean some of them are not clean enough .
- u_{11} that 's my greatest concern . but how can you tell which one is clean ?
- u_{12} i judge by appearances .
- u_{13} i got it . i think it 's worth trying .
- u_{14} it certainly is .
- u_{15} i 'd like to try some kebab , roasted squid , and many different appealing things .
- u_{16} do n't try everything at one time , please .
- u_{17} i see . thank you .

A sample dialog from DailyDialog (Li et al., 2017) of length 17.

Experiments and Results

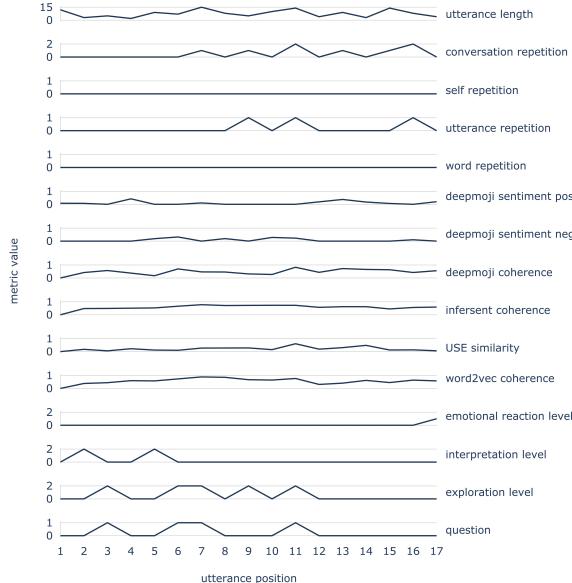
Analysis of metric dynamics in written human conversations



Maiborn
Wolff
Mensch IT

- u_1 did you ever buy food from the snack stands near our hotel ?
 u_2 yes , several times .
 u_3 how do you like them ?
 u_4 not bad .
 u_5 i always have the temptation to eat something there .
 u_6 then , why did n't you do that ?
 u_7 i do n't know how much we can trust them . do you have any ideas ?
 u_8 some of them , i think , are not good .
 u_9 it does n't taste good ?
 u_{10} no , i mean some of them are not clean enough .
 u_{11} that 's my greatest concern . but how can you tell which one is clean ?
 u_{12} i judge by appearances .
 u_{13} i got it . i think it 's worth trying .
 u_{14} it certainly is .
 u_{15} i 'd like to try some kebab , roasted squid , and many different appealing things .
 u_{16} do n't try everything at one time , please .
 u_{17} i see . thank you .

A sample dialog from DailyDialog (Li et al., 2017) of length 17.



The metric progress within the sample dialog.

Experiments and Results

Analysis of metric dynamics in written human conversations



Dialog 1

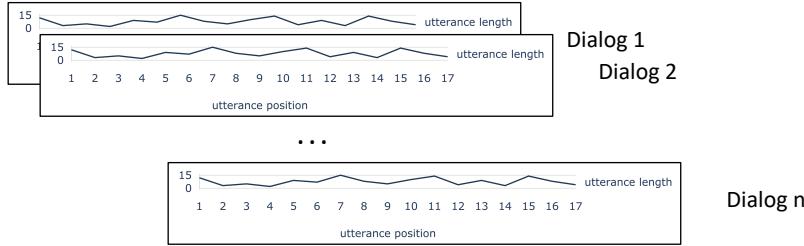
σ : Intra-dialog metric standard deviation

μ : Metric mean across the positions in the dialog

Coefficient of Variation (CV): $CV = \sigma / \mu$

Experiments and Results

Analysis of metric dynamics in written human conversations



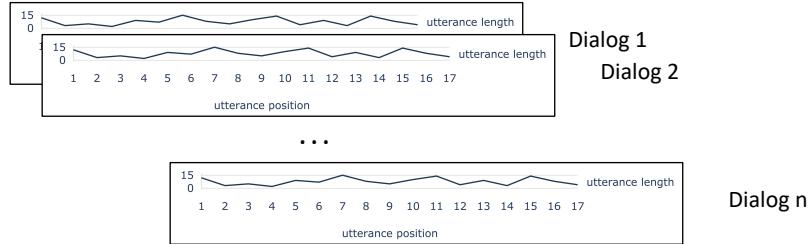
σ : Average intra-dialog metric standard deviation

μ : Average metric mean across the positions in the dialog

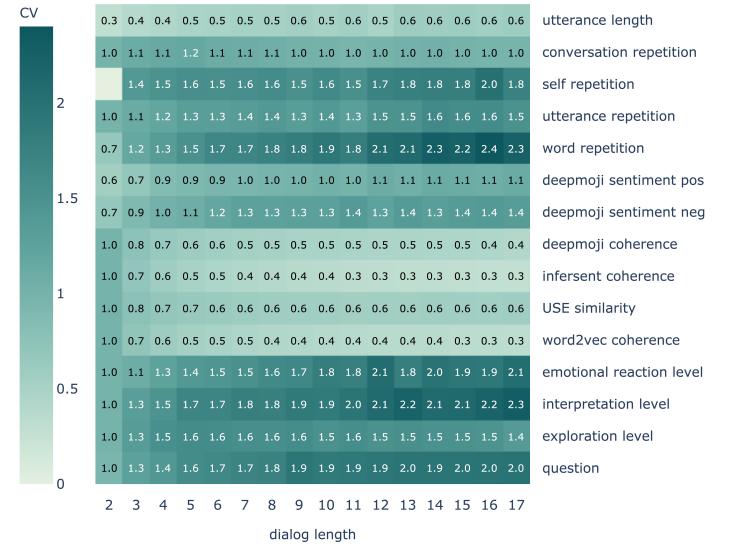
Coefficient of Variation (CV): $CV_{l,m} = \sigma / \mu$

Experiments and Results

Analysis of metric dynamics in written human conversations



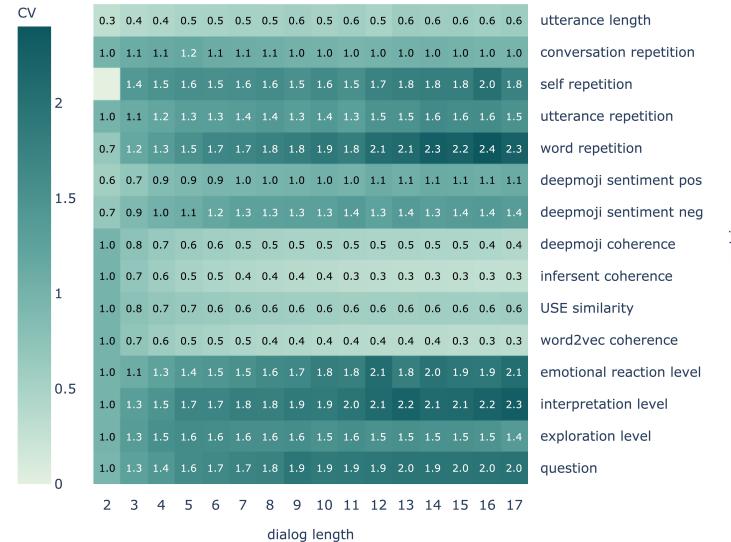
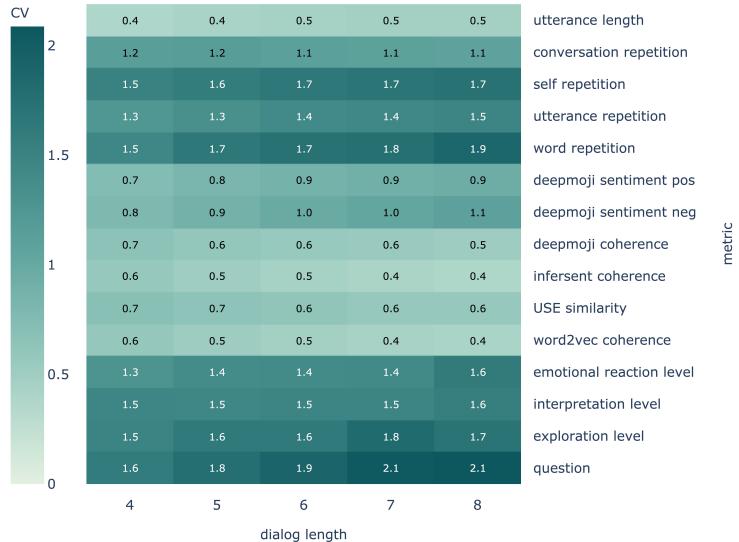
σ : Average intra-dialog metric standard deviation
 μ : Average metric mean across the positions in the dialog
 Coefficient of Variation (CV): $CV_{l,m} = \sigma / \mu$



Heatmap of $CV_{l,m}$, the average variability of metric m within all dialogs of length l , for all dialog lengths and all metrics in the dialogs in DailyDialog.

Experiments and Results

Analysis of metric dynamics in written human conversations



Experiments and Results

1. Analysis of metric dynamics in written human conversations → Metrics in human dialogs are dynamic

2. **DYME**: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Experiments and Results

1. Analysis of metric dynamics in written human conversations → Metrics in human dialogs are dynamic

2. **DYME**: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

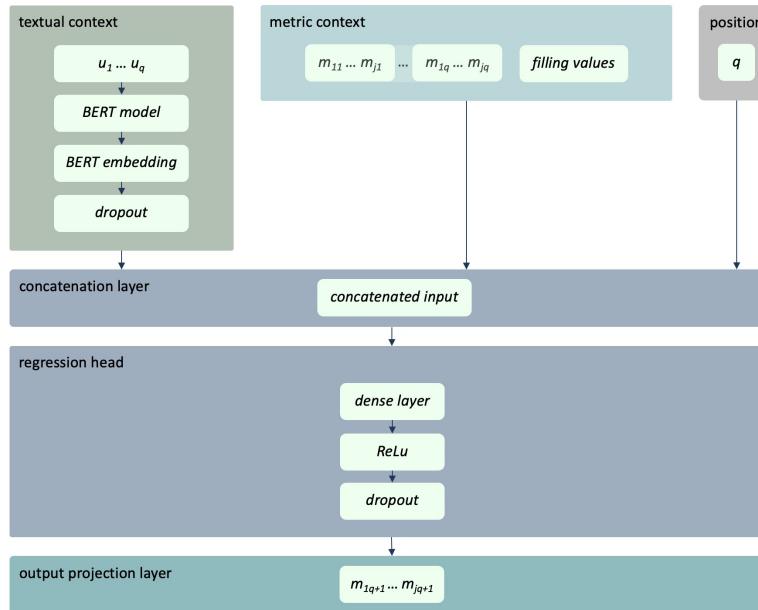
Task

- Given a dialog history (textual and metric context) predict the values of all metrics for the next utterance
- Multi-label regression

Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Approach: DYME – A dynamic metric

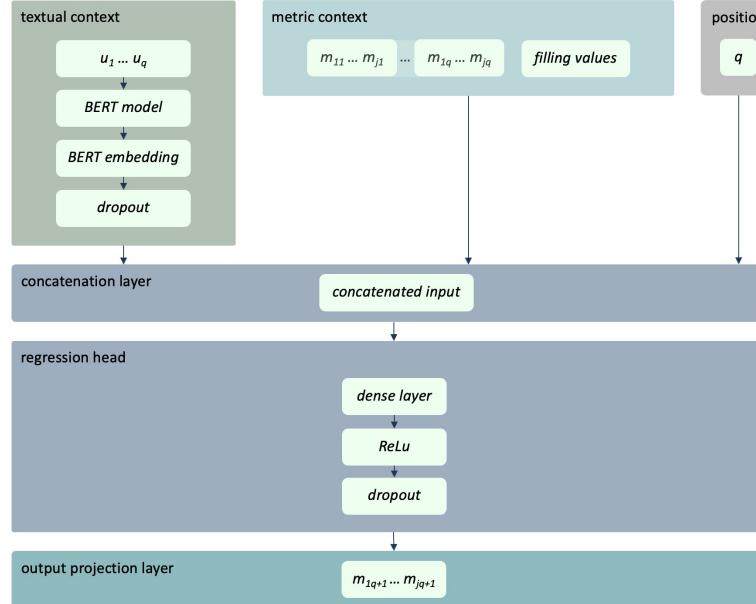


Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Approach: DYME – A dynamic metric

u_1 did you ever buy food from the snack stands near our hotel ?
 u_2 yes , several times .
 u_3 how do you like them ?
 u_4 not bad .
 u_5 i always have the temptation to eat something there .
 u_6 then , why did n't you do that ?
 u_7 i do n't know how much we can trust them . do you have any ideas ?
 u_8 some of them , i think , are not good .

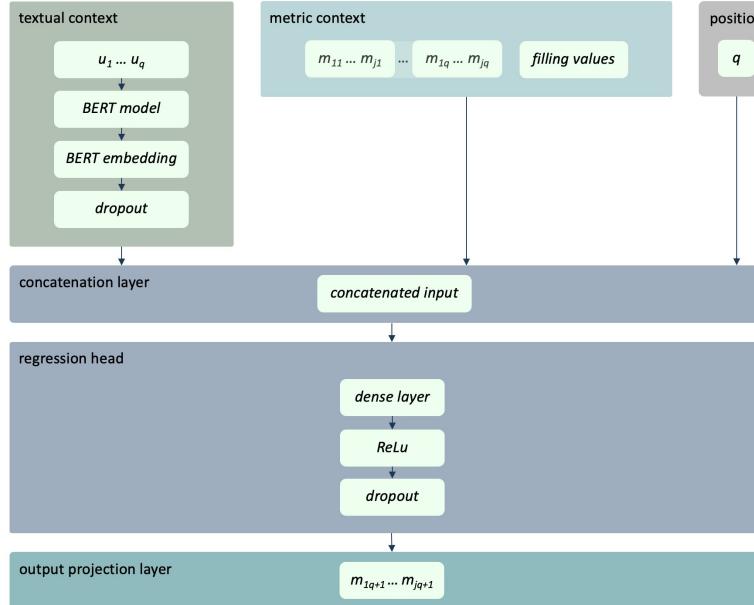


Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Approach: DYME – A dynamic metric

u_1 did you ever buy food from the snack stands near our hotel ?
 u_2 yes , several times .
 u_3 how do you like them ?
 u_4 not bad .
 u_5 i always have the temptation to eat something there .
 u_6 then , why did n't you do that ?
 u_7 i do n't know how much we can trust them . do you have any ideas ?
 u_8 some of them , i think , are not good .



Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

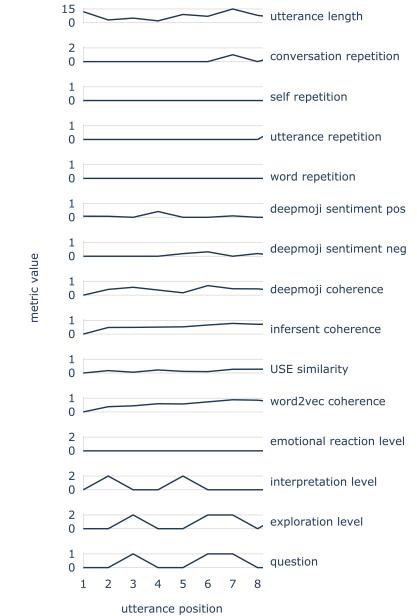
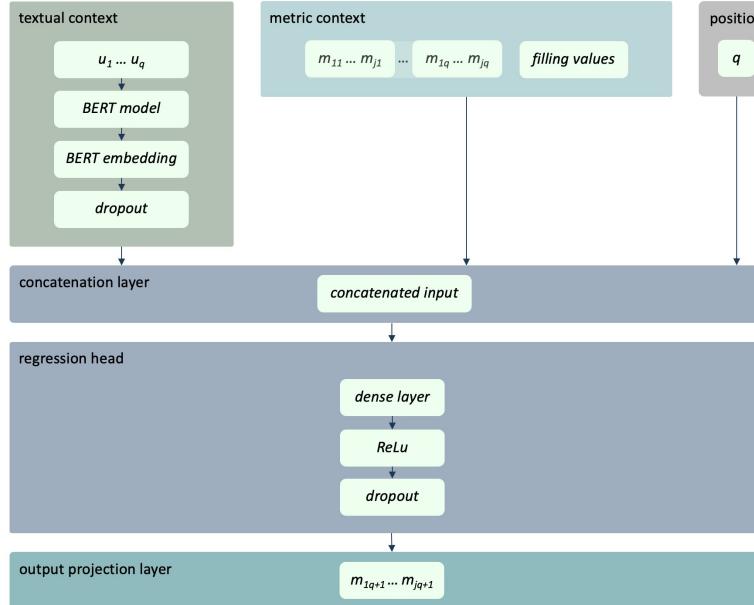


Maiborn
Wolff

MenschIT

Approach: DYME – A dynamic metric

u_1 did you ever buy food from the snack stands near our hotel ?
 u_2 yes , several times .
 u_3 how do you like them ?
 u_4 not bad .
 u_5 i always have the temptation to eat something there .
 u_6 then , why did n't you do that ?
 u_7 i do n't know how much we can trust them . do you have any ideas ?
 u_8 some of them , i think , are not good .



Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

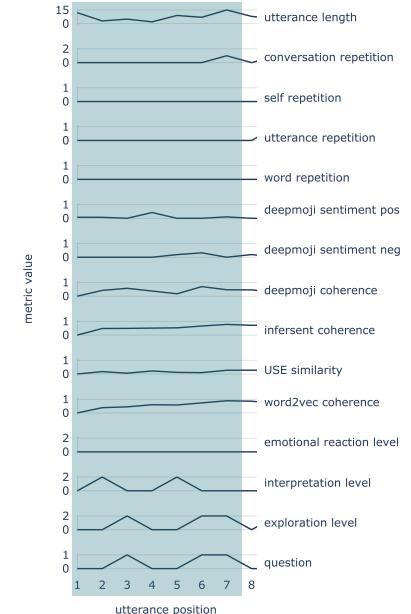
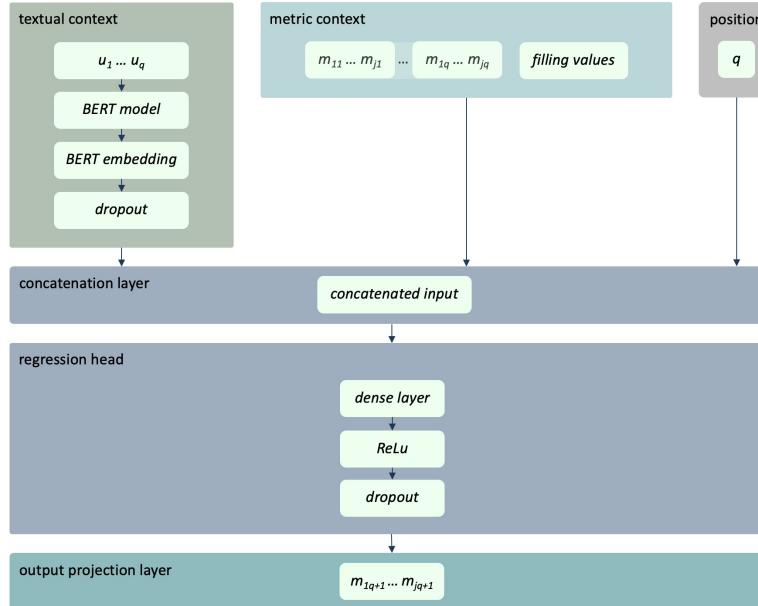


Maiborn
Wolff

MenschIT

Approach: DYME – A dynamic metric

u_1 did you ever buy food from the snack stands near our hotel ?
 u_2 yes , several times .
 u_3 how do you like them ?
 u_4 not bad .
 u_5 i always have the temptation to eat something there .
 u_6 then , why did n't you do that ?
 u_7 i do n't know how much we can trust them . do you have any ideas ?
 u_8 some of them , i think , are not good .



Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

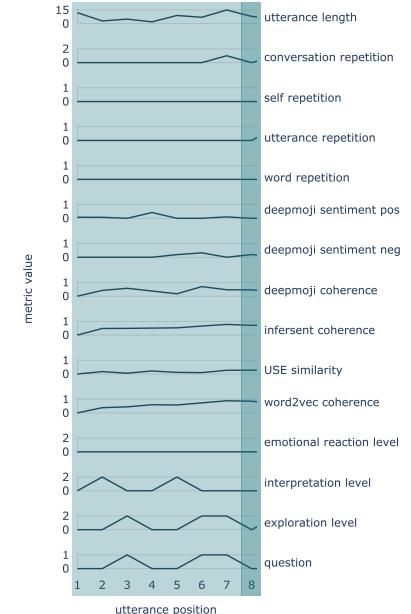
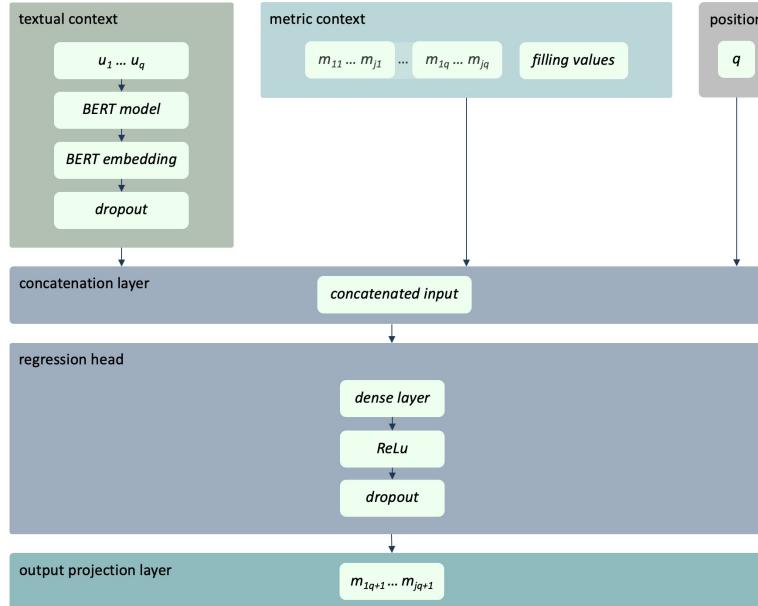


Maiborn
Wolff

MenschIT

Approach: DYME – A dynamic metric

u_1 did you ever buy food from the snack stands near our hotel ?
 u_2 yes , several times .
 u_3 how do you like them ?
 u_4 not bad .
 u_5 i always have the temptation to eat something there .
 u_6 then , why did n't you do that ?
 u_7 i do n't know how much we can trust them . do you have any ideas ?
 u_8 some of them , i think , are not good .



Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Training

- Supervised Learning
- DailyDialog and EmpatheticDialogues (80% training, 10% validation, 10% testing)
- Learning objective: minimize the average Mean Squared Error (MSE) between the predicted metric values \hat{y} and the ground truth metric values y across all metrics

$$loss(y, \hat{y}) = \frac{1}{j} \sum_{i=1}^j (y_i - \hat{y}_i)^2$$

Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Baseline

- Moving average baseline
- Uses the per-metric average of the given metric context as prediction for the next utterance position
- Does not make use of the textual context at all
- **Would achieve perfect accuracy if metrics were static**

Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Results (Average RMSE on test dataset)

Baseline: 0.2326

DYME: 0.1845

→ Absolute improvement: 4.81%

→ Relative improvement: 20.68%

Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Results (Average RMSE on test dataset)

Baseline: 0.2326
DYME: 0.1845

→ Absolute improvement: 4.81%
→ Relative improvement: 20.68%

Metric	DYME	Baseline	Absolute improvement over baseline	Relative improvement over baseline (in %)
Word2Vec coherence	0.1237	0.2575	0.1338	51,96
Inferent coherence	0.1086	0.2232	0.1146	51,34
Deepemoji coherence	0.2114	0.2879	0.0765	26,57
Question	0.2121	0.2804	0.0683	24,36
USE similarity	0.1552	0.2002	0.0450	22,48
Exploration level	0.3048	0.3599	0.0551	15,31
Emotional reaction level	0.2621	0.3059	0.0438	14,32
Interpretation level	0.3370	0.3791	0.0421	11,11
Word repetition	0.0600	0.0671	0.0071	10,58
Conversation repetition	0.0781	0.0854	0.0073	8,55
Deepemoji sentiment	0.1568	0.1700	0.0041	7,76
Self repetition	0.1068	0.1149	0.0081	7,05
Utterance repetition	0.0954	0.1020	0.0066	6,47
Utterance length	0.1039	0.0954	-0.0085	-8,91

Average prediction errors (average RMSE) across all positions for each metric in descending order of the improvement of DYME over the baseline. The reported value for a metric m describes the model's average prediction error across all positions when predicting the value of metric m.

Experiments and Results

DYME: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations

Results (Average RMSE on test dataset)

Baseline: 0.2326

DYME: 0.1845

→ Absolute improvement: 4.81%

→ Relative improvement: 20.68%

Position	DYME	Baseline	Absolute improvement over baseline	Relative improvement over baseline (in %)
1	0.1523	0.3192	0.1669	52.29
2	0.1619	0.2490	0.0871	34.98
3	0.1905	0.2381	0.0476	19.99
4	0.1792	0.2223	0.0431	19.39
5	0.1703	0.2081	0.0378	18.16
6	0.1792	0.2049	0.0257	12.54
7	0.1700	0.1992	0.0292	14.66

Average prediction errors (average RMSE) across all metrics at each prediction position in ascending order of the position. The reported value for a prediction position q describes the model's average prediction error across all metrics when predicting for utterance position $q+1$, given the input context up to utterance position q .

Conclusion

1. Metrics in human dialogs are dynamic.
2. A dynamic metric can be learned from dialog data.

We propose **DYME**: A DYnamic MEtric for Dialog Modeling Learned from Human Conversations.

References

- Brown, C.E.: Coefficient of variation. In: Applied multivariate statistics in geohydrology and related sciences, pp. 155–157. Springer (1998)
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: 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). pp. 4171–4186 (2019)
- Li, J., Monroe, W., Ritter, A., Jurafsky, D., Galley, M., Gao, J.: Deep reinforcement learning for dialogue generation. In: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing. pp. 1192– 1202. Association for Computational Linguistics, Austin, Texas (Nov 2016). <https://doi.org/10.18653/v1/D16-1127>
- Li, Y., Su, H., Shen, X., Li, W., Cao, Z., Niu, S.: DailyDialog: A manually labelled multi-turn dialogue dataset. In: Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 986–995 (2017)
- Rashkin, H., Smith, E.M., Li, M., Bureau, Y.L.: Towards empathetic open-domain conversation models: A new benchmark and dataset. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. pp. 5370–5381 (2019)
- Saleh, A., Jaques, N., Ghandeharioun, A., Shen, J., Picard, R.: Hierarchical reinforcement learning for open-domain dialog. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 8741–8748 (2020)
- Sharma, A., Miner, A., Atkins, D., Althoff, T.: A computational approach to understanding empathy expressed in text-based mental health support. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 5263–5276 (2020)
- Xu, C., Wu, W., Wu, Y.: Towards explainable and controllable open domain dialogue generation with dialogue acts. arXiv preprint arXiv:1807.07255 (2018)

Thank you!

› Slides, code and models on GitHub: <https://github.com/florianvonunold/DYME>

