Microsoft Research AI

# Dialogue Response Ranking Training with Large-Scale Human Feedback Data

**Xiang Gao**, Yizhe Zhang, Michel Galley, Chris Brockett, Bill Dolan
*Microsoft Research AI, Redmond, WA, USA*

**EMNLP 2020**

The 2020 Conference on Empirical Methods
in Natural Language Processing

1

- Great progress in building conversational AI with large-scale pre-trained models
- They are trained mostly by minimizing **perplexity** on human samples

Microsoft

**DialoGPT**: arxiv.org/abs/1911.00536
Trained with **147 M** Reddit Dialogues!

Google Brain

**Meena**: arxiv.org/abs/2001.09977

FACEBOOK

**Blender**: arxiv.org/abs/2004.13637

*And many other awesome works..*

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

**2**

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

**EMNLP 2020**
The 2020 Conference on Empirical Methods in Natural Language Processing
*16th – 20th November 2020*

# Motivation

- Great progress in building conversational AI with large-scale pre-trained models

- They are trained mostly by minimizing **perplexity** on human samples

- However, some human replies are more engaging than others, spawning more follow-up interactions
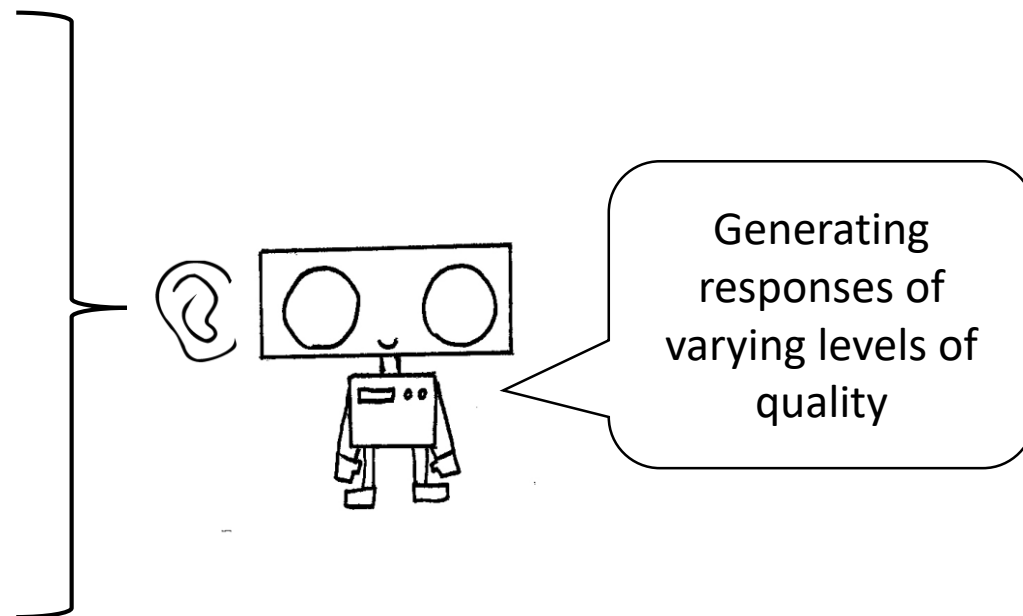


engaging

Boring/bland/generic

Hate/offensive/toxic

Generating responses of varying levels of quality

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

**EMNLP 2020**
The 2020 Conference on Empirical Methods in Natural Language Processing
*16th – 20th November 2020*

# Motivation

- Existing ranking methods may be suboptimal

  - **Perplexity**: e.g. MMI, only reflects relevancy
  - **Manually designed features**: not directly based on real-world human preferences in an end-to-end fashion.

- Crowdsourcing of large-scale training data is too expensive

- Social networks provide ways to measure Human feedback on dialogues (and other contents).



**A Diversity-Promoting Objective Function for Neural Conversation Models**

Jiwei Li[1*]   Michel Galley[2]   Chris Brockett[2]   Jianfeng Gao[2]   Bill Dolan[2]

[1]Stanford University, Stanford, CA, USA
jiweil@stanford.edu
[2]Microsoft Research, Redmond, WA, USA
{mgalley,chrisbkt,jfgao,billdol}@microsoft.com

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

**4**

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

EMNLP 2020
The 2020 Conference on Empirical Methods in Natural Language Processing
16th – 20th November 2020

- Optimizing expected **_human feedback_**, not just _perplexity_?
- Social network human feedback data!

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

**5**

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

EMNLP 2020
The 2020 Conference on Empirical Methods in Natural Language Processing
16th – 20th November 2020

- **Motivation**
- [**Dataset**](#)
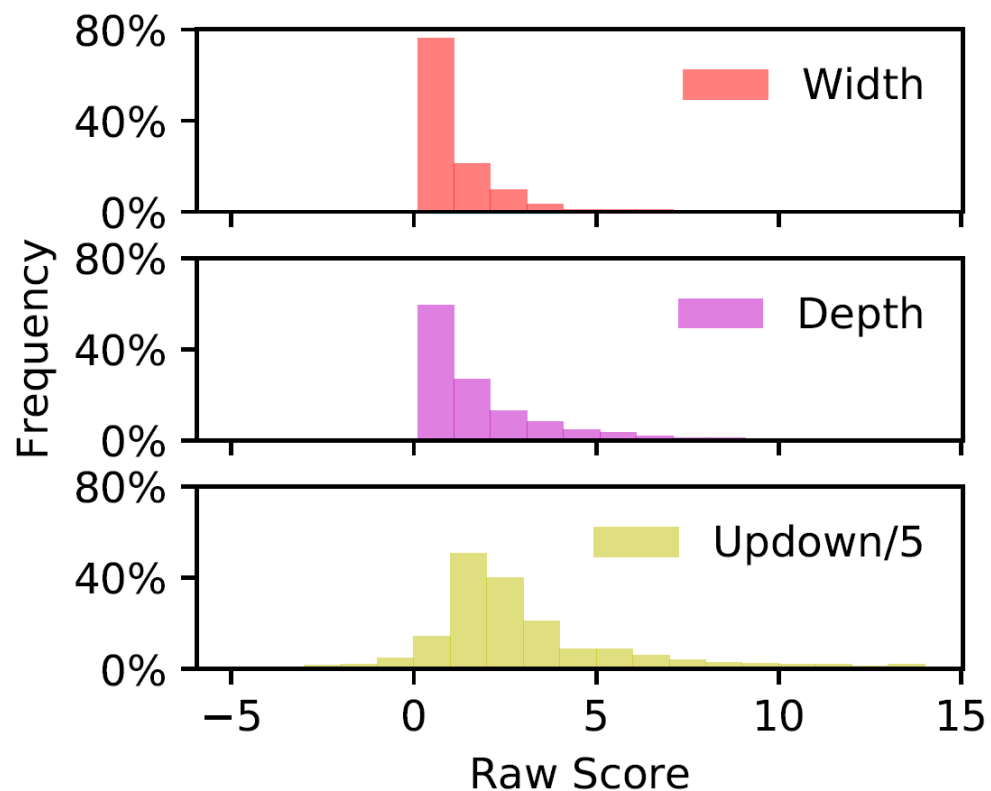- **Method**
- **Results**

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**
Xiang Gao, Yizhe Zhang, Michel Galley, Chris Brockett, Bill Dolan
*Microsoft Research AI, Redmond, WA, USA*

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

EMNLP 2020

The 2020 Conference on Empirical Methods
in Natural Language Processing

6

*16th – 20th November 2020*

# Dataset

- We define three metrics of human feedback on Reddit
  - Updown
  - Width
  - Depth

$u_0$: I love NLP!
👍20 👎2

"Updown" is 17 – 3 = 14

$u_1$: Here's a great NLP textbook (URL)
👍17 👎3

$u_2$: Me too!
👍2 👎0

$u_3$: Anything focused on dialog?
👍3 👎0

$u_4$: Thanks!
👍0 👎0

$u_5$: Awesome!
👍0 👎0

"Depth" of $u_1$ is 2

"Width" of $u_1$ is 3

$u_6$: Sure, here you are (URL)
👍12 👎0

- All of three metrics have a long-tailed distribution

- They are correlated. Width and depth are more correlated as both are measure of number of replies



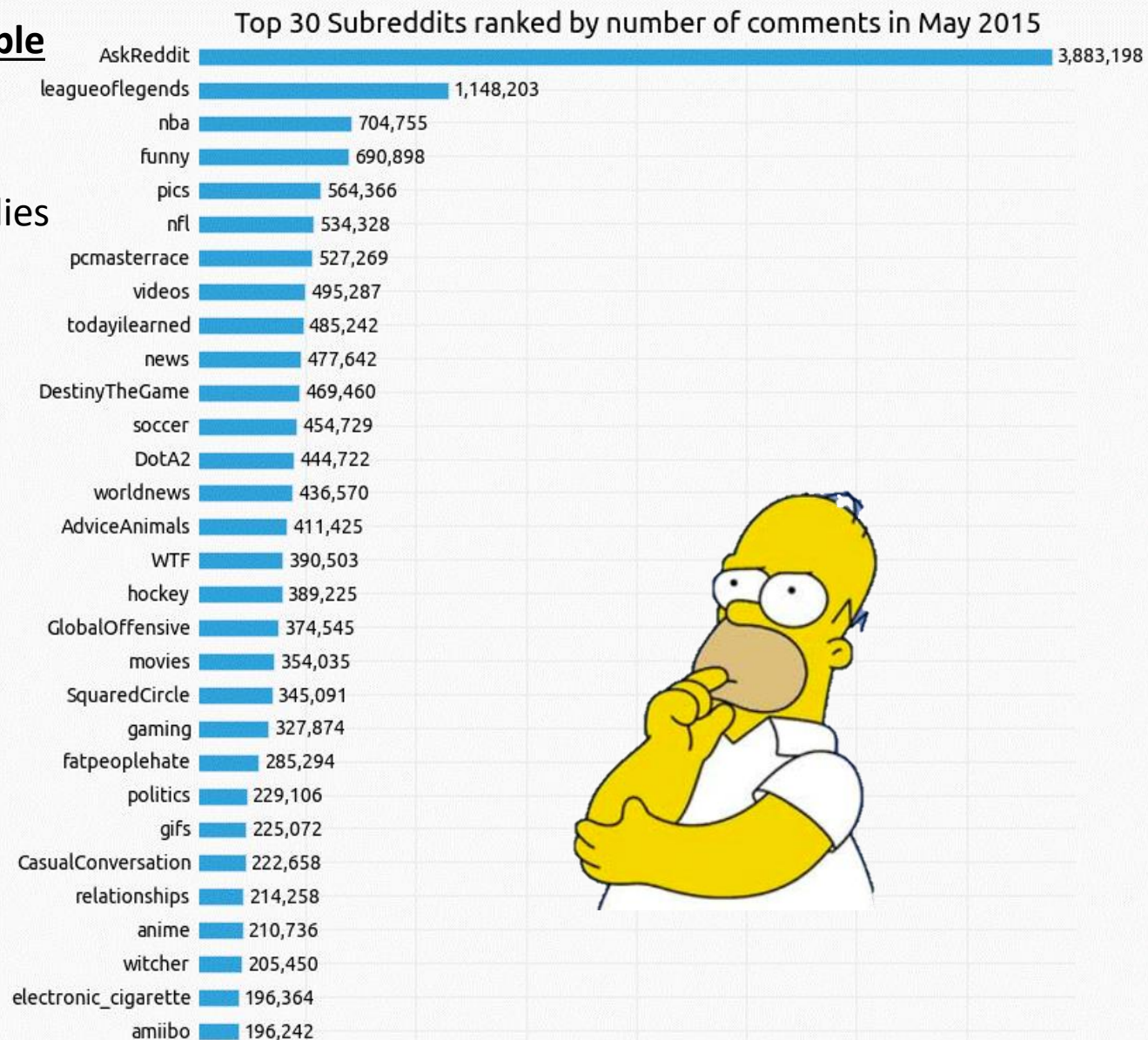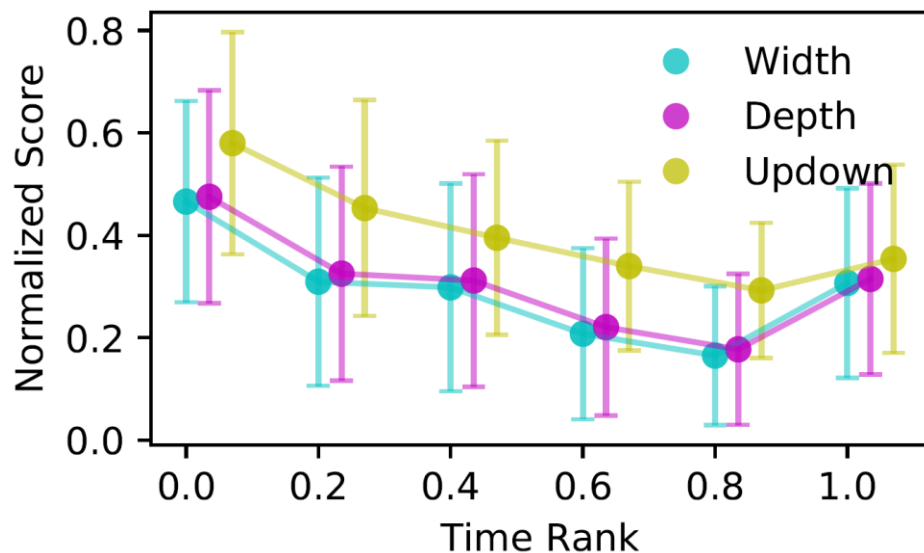| | Width | Depth | Updown |
|---|---|---|---|
| Width | 1 | 0.8592 | 0.3491 |
| Depth | 0.8592 | 1 | 0.3257 |
| Updown | 0.3491 | 0.3257 | 1 |

Table 1: Spearman's $\rho$ between different measurements of human feedback. Darker cell color indicates higher correlation.

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

8

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

**EMNLP 2020**
The 2020 Conference on Empirical Methods in Natural Language Processing
*16th – 20th November 2020*

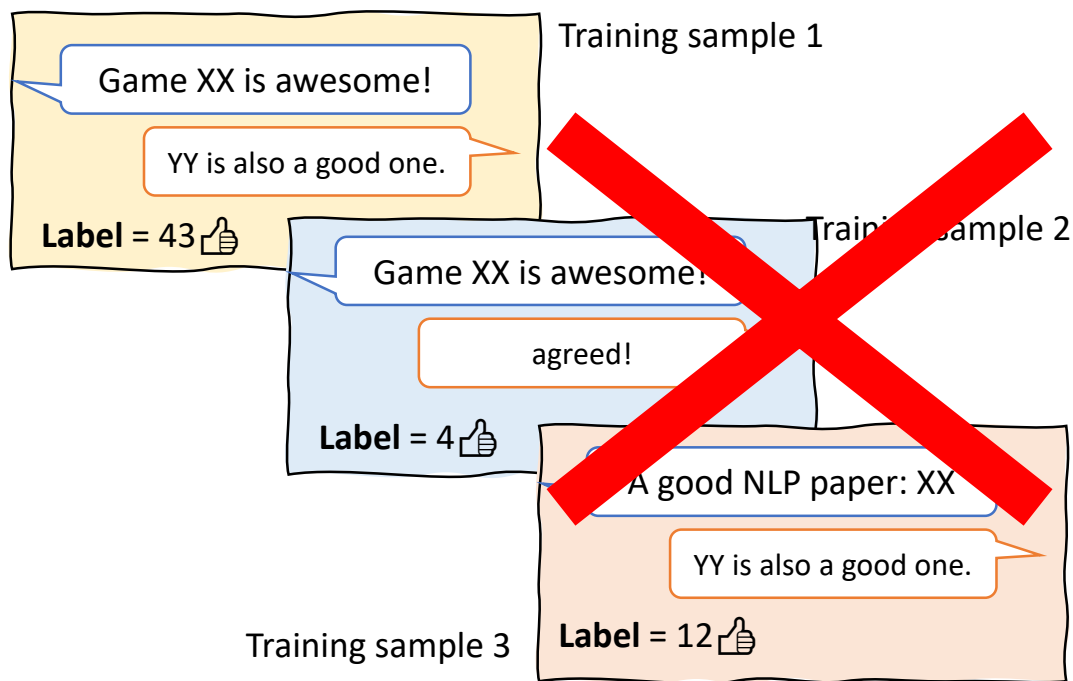- However, these metrics are **not directly usable/comparable**
- Confounding factors
  - **Topics/Subreddits**: popularity differs significantly
  - **Timing**: The early comments gets more upvotes/replies





Top 30 Subreddits ranked by number of comments in May 2015

| Subreddit | Comments |
|---|---|
| AskReddit | 3,883,198 |
| leagueoflegends | 1,148,203 |
| nba | 704,755 |
| funny | 690,898 |
| pics | 564,366 |
| nfl | 534,328 |
| pcmasterrace | 527,269 |
| videos | 495,287 |
| todayilearned | 485,242 |
| news | 477,642 |
| DestinyTheGame | 469,460 |
| soccer | 454,729 |
| DotA2 | 444,722 |
| worldnews | 436,570 |
| AdviceAnimals | 411,425 |
| WTF | 390,503 |
| hockey | 389,225 |
| GlobalOffensive | 374,545 |
| movies | 354,035 |
| SquaredCircle | 345,091 |
| gaming | 327,874 |
| fatpeoplehate | 285,294 |
| politics | 229,106 |
| gifs | 225,072 |
| CasualConversation | 222,658 |
| relationships | 214,258 |
| anime | 210,736 |
| witcher | 205,450 |
| electronic_cigarette | 196,364 |
| amiibo | 196,242 |

Author: Ramiro Gómez - ramiro.org • Data: Reddit /u/Stuck_In_the_Matrix & /u/fhoffa - reddit.com

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

- **Motivation**
- **Dataset**
- **Method**
- **Results**

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**
Xiang Gao, Yizhe Zhang, Michel Galley, Chris Brockett, Bill Dolan
*Microsoft Research AI, Redmond, WA, USA*

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

**EMNLP 2020**

The 2020 Conference on Empirical Methods
in Natural Language Processing

# Contrastive Dataset

- Directly predicting metric value is **hard**, due to confounding factors (e.g. timing of post)

- **Contrastive learning!**



**Predicting feedback metric value**

Training sample 1

Game XX is awesome!

YY is also a good one.

**Label** = 43 👍

Training sample 2

Game XX is awesome!

agreed!

**Label** = 4 👍

A good NLP paper: XX

YY is also a good one.

Training sample 3    **Label** = 12 👍

**Classify which one gets more feedback**

[**Context**] Game XX is awesome!

[**Hyp-A**]: YY is also a good one.    [**Hyp-B**]: Agreed!

43 👍    4 👍

**Label** = hyp-A is better

**Only compare pairs of responses that are comparable**
- For the same dialogue context
- Published at roughly the same time
- ...

# Contrastive Learning

$h$ = **Logits** (scalar)

$c$ = **Context** (string)

$$h(c, r) = \text{DialogRPT}(c, r)$$

$r$ = **Response** (string)

- Inferred Score

$$s(r|c) = \text{Sigmoid}(h(c, r))$$

- Training Loss

$$\mathcal{L} = -\sum_{i \in \text{batch}} \log \frac{e^{h(c_i, r_i^+)}}{e^{h(c_i, r_i^+)} + e^{h(c_i, r_i^-)}}$$

negative log likelihood

$= P(r^+|c)$
Softmax probability to pick $r^+$
- $r^+$ for positive sample
- $r^-$ for negative sample

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

# Implementation

- GPT-2 type model, initialized with DialoGPT weight
- Using the latent vector at the last time step to compute logit and score
- **DialogRPT**: <u>D</u>ialogue <u>R</u>anking <u>P</u>retrained <u>T</u>ransformers



Logits $h(c, r)$

linear

Latent vector

Transformer-12

...

Transformer-2

Transformer-1

Position embed

+

0   1   2   3   4   5   6   7

Token embed

*I*   *love*   *NLP*   *!*   *<endoftext>*   *Me*   *too*   *<endoftext>*

Context $c$          Response $r$

Dialogue Response Ranking Training with
Large-Scale Human Feedback Data

EMNLP 2020
The 2020 Conference on Empirical Methods
in Natural Language Processing
16th – 20th November 2020

**However, can we apply rankers trained on human vs human data on generators?**

$r$ is a human-like response: $\quad r \in H$

$\quad$ =0, assumed

Probability that response $r$ gets the most feedback given context $c$:

$$P(r|c) = P(r|c, r \in H)P(r \in H) + P(r|c, r \notin H)P(r \notin H)$$

$$= P(r|c, r \in H)P(r \in H)$$

| Task | Subtask Description | Training size (number of pairs) |
|---|---|---|
| $P(r \in H)$ **Human vs fake** | Fake = Retrieved human response | 40.7 M |
| | Fake = Machine generated response | 40.7 M |
| $P(r|c, r \in H)$ **Human vs human** (which gets more feedback) | Feedback = Updown (more upvotes - downvotes) | 40.7 M |
| | Feedback = Width (more direct replies) | 22.3 M |
| | Feedback = Depth (longer follow-up thread) | 25.1 M |

14

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

- **Motivation**
- **Dataset**
- **Method**
- **Results**

EMNLP 2020

The 2020 Conference on Empirical Methods
in Natural Language Processing

15

*16th – 20th November 2020*

| Context: I love NLP! | | | |
|---|---|---|---|
| Response: | Width | Depth | Updown |
| A    Me too! | 0.033 | 0.043 | 0.171 |
| B    It's super useful and more and more powerful! | 0.054 | 0.164 | 0.296 |
| C    Can you tell me how it works? | 0.644 | **0.696** | 0.348 |
| D    Can anyone recommend a nice review paper? | **0.687** | 0.562 | 0.332 |
| E    Here's a free textbook (URL) in case anyone needs it. | 0.319 | 0.409 | **0.612** |

Table 3: Predicted feedback scores of several example responses given the same context.

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

**16**

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

**EMNLP 2020**
The 2020 Conference on Empirical Methods in Natural Language Processing
*16th – 20th November 2020*

# Generator reranking

Although hypothesis C is most likely to be generated (Generation Probability = 0.496), it's relatively boring.
Using Updown Score, we can pick the hypothesis A, which is perhaps more interesting (Updown Score = 0.431)

**[Context]**: Can we restart 2020?

| | Generation Probability | Updown Score | Generated Hypothesis |
|---|---|---|---|
| A | 0.383 | 0.431 | I think we should go back to the beginning, and start from the beginning. |
| B | 0.195 | 0.323 | I think we should just give up and let the year just pass. |
| C | 0.496 | 0.302 | Yes, we can. |
| D | 0.328 | 0.153 | I think so, yes. |

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

17

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

**EMNLP 2020**
The 2020 Conference on Empirical Methods in Natural Language Processing
*16th – 20th November 2020*

| | Method | Pairwise accuracy | Spearman $\rho$ |
|---|---|---|---|
| Width | Dialog ppl. | 0.513 | -0.009 |
| | Reverse dialog ppl. | 0.571 | 0.099 |
| | Length baseline | 0.595 | 0.229 |
| | BoW baseline | 0.596 | 0.234 |
| | DIALOGRPT | **0.752** | **0.357** |
| Depth | Dialog ppl. | 0.508 | -0.004 |
| | Reverse dialog ppl. | 0.557 | 0.063 |
| | Length baseline | 0.543 | 0.134 |
| | BoW baseline | 0.584 | 0.187 |
| | DIALOGRPT | **0.695** | **0.317** |
| Updown | Dialog ppl. | 0.488 | 0.003 |
| | Reverse dialog ppl. | 0.560 | 0.076 |
| | Length baseline | 0.531 | 0.063 |
| | BoW baseline | 0.571 | 0.134 |
| | DIALOGRPT | **0.683** | **0.295** |

Table 5: Performance on test set ranking gold responses, measured by pairwise accuracy and Spearman's $\rho$.

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

18

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

**EMNLP 2020**
The 2020 Conference on Empirical Methods in Natural Language Processing
16th – 20th November 2020

For each context, there're $k$ human responses and $n$ distractor (random human responses),
Rank these $k + n$ candidates based on predicted $P(r \in H)$
Even DialogRPT is only trained on Reddit, it performs very well on all four datasets

| Dataset | Method | Hits@$k$ |
|---|---|---|
| Reddit ($k > 5, n=k$) | BLEU1 | 0.651 |
| | BERTScore | 0.685 |
| | BLEURT | 0.714 |
| | BM25 | 0.309 |
| | ConvRT | 0.760 |
| | Dialog ppl. | 0.560 |
| | Reverse dialog ppl. | 0.775 |
| | DIALOGRPT | **0.886** |
| DailyDialog ($k=1, n=19$) | BM25 | 0.182 |
| | ConvRT | 0.380 |
| | Dialog ppl. | 0.176 |
| | Reverse dialog ppl. | 0.457 |
| | DIALOGRPT | **0.621** |

Zero-short!

| Dataset | Method | Hits@$k$ |
|---|---|---|
| Twitter ($k=1, n=19$) | BM25 | 0.178 |
| | ConvRT | 0.439 |
| | Dialog ppl. | 0.107 |
| | Reverse dialog ppl. | 0.440 |
| | DIALOGRPT | **0.548** |
| PersonaChat ($k=1, n=19$) | BM25 | 0.117 |
| | ConvRT | 0.197 |
| | IR Baseline | 0.213 |
| | Starspace | 0.318 |
| | KV profile memory | 0.349 |
| | Dialog ppl. | 0.108 |
| | Reverse dialog ppl. | 0.449 |
| | DIALOGRPT | **0.479** |

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

**19**

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

**EMNLP 2020**
The 2020 Conference on Empirical Methods in Natural Language Processing
*16th – 20th November 2020*

| Model | Trained on | Tested on | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | Human vs. Human | | | Human vs. Fake | |
| | | Width | Depth | Updown | Rand | Generated |
| Human feedback | Width | 0.764 | 0.693 | 0.601 | 0.517 | 0.644 |
| | Depth | 0.749 | 0.701 | 0.588 | 0.512 | 0.647 |
| | Updown | 0.659 | 0.602 | 0.683 | 0.526 | 0.667 |
| Human-like | Rand | 0.558 | 0.552 | 0.522 | 0.843 | 0.413 |
| | + Generated | 0.560 | 0.558 | 0.522 | 0.864 | 0.880 |
| Ensemble | - | 0.746 | 0.675 | 0.666 | 0.758 | 0.821 |

Table 6: Pairwise accuracy of DIALOGRPT models. Darker cell color indicates better performance.

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io

20

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

EMNLP 2020
The 2020 Conference on Empirical Methods in Natural Language Processing
*16th – 20th November 2020*

# Open-sourced!

Dataset available at:
https://dialogfeedback.github.io

Pre-trained models at:
https://github.com/golsun/DialogRPT

# Open-sourced!



Demo available at:
http://github.com/golsun/DialogRPT

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**

**EMNLP 2020**
The 2020 Conference on Empirical Methods in Natural Language Processing
*16th – 20th November 2020*

# Thank you!

**Dialogue Response Ranking Training with Large-Scale Human Feedback Data**
Xiang Gao, Yizhe Zhang, Michel Galley, Chris Brockett, Bill Dolan
*Microsoft Research AI, Redmond, WA, USA*

paper: arxiv.org/abs/2009.06978
code: github.com/golsun/**DialogRPT**
data: https://dialogfeedback.github.io