# Dialogue Generation on Infrequent Sentence Functions via Structured Meta-Learning
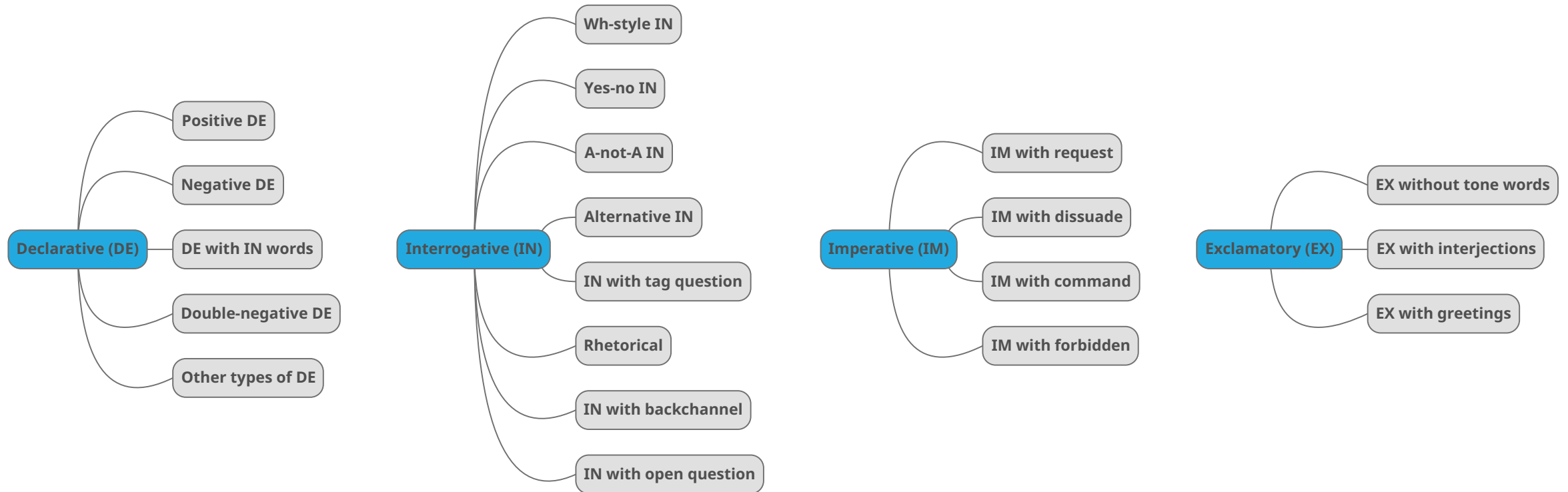
**Yifan Gao**[1], Piji Li[2], Wei Bi[2], Xiaojiang Liu[2], Michael R. Lyu[1], and Irwin King[1]

1. The Chinese University of Hong Kong     2. Tencent AI Lab

# Introduction

Sentence Function

Definition: "In linguistics, a sentence function refers to a speaker's purpose in uttering a specific sentence, phrase, or clause." [Wikipedia]

Findings of ACL: EMNLP 2020, Fine-Grained Sentence Functions for Short-Text Conversation

# Introduction

## Sentence Functions in Conversation

- Humans express intentions in conversations through sentence functions.

- Sentence functions have great influences on the structures of utterances in conversations including word orders, syntactic patterns, and other aspects.

| Sentence Function | Frequent Patterns | | Sentence Examples | |
|---|---|---|---|---|
| | Chinese | English | Chinese | English |
| Wh-style IN | x在哪y?<br>谁是x? | Where does x y?<br>Who is x? | 周末在哪过啊<br>谁是天蝎座 | <u>Where</u> do you spend your weekend<br><u>Who</u> is a Scorpio |
| Yes-no IN | x是在y吗?<br>x是指y吗? | Is x y?<br>Does x y? | 你是在云南吗?<br>你是指昨天的篮球比赛吗? | <u>Are</u> you in Yunnan?<br><u>Do</u> you <u>mean</u> the basketball match yesterday? |
| Alternative IN | x还是y<br>x y哪个 | x or y<br>x y which | 狮子和白羊真配还是假配?<br>香蕉和苹果哪个卖得比较好? | Leo and Aries go together <u>or</u> not?<br><u>Which</u> sells better, banana or apple? |

Frequent word patterns of three level-2 Interrogative sentence functions. (Bi et al., ACL 2019)

# Introduction

## Imbalance Problem in Large-scale Conversation Dataset with Sentence Function Annotation

- Existing work shows that the use of sentence functions improves the overall quality of generated responses (Ke et al., ACL 2018).

- However, the number of utterances for different types of fine-grained sentence functions is usually extremely imbalanced. In the large-scale dataset STC-SeFun (Bi et al., ACL 2019):

| Sentence Function | Query | Response |
|---|---|---|
| Declarative (DE) | | |
| Positive DE | 49,223 (48%) | 67,540 (57%) |
| Negative DE | 9,241(9%) | 18,428(16%) |
| DE with IN words | 887(.9%) | 2,660(2%) |
| Double-negative DE | 40(<.1%) | 99(.1%) |
| Other types of DE | 2,675(3%) | 5,218(4%) |

**Dialogue generation models suffer from data deficiency for these infrequent sentence functions!**

Findings of ACL: EMNLP 2020, Fine-Grained Sentence Functions for Short-Text Conversation
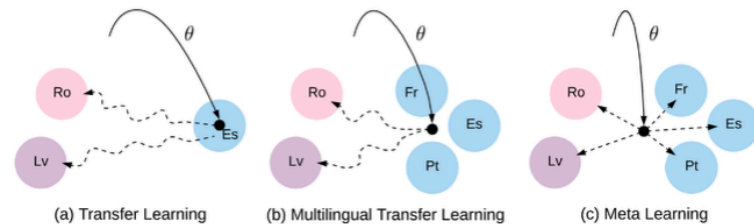
# **Proposed Approach**

## Model-Agnostic Meta-Learning (MAML)

Model-Agnostic Meta-Learning (MAML) can learn from a variety of **tasks** such that it can solve new learning tasks using only **a small number** of training samples.

- Training: Learn transferable internal representations across tasks (task=domain).

- Testing: Quickly adapt to a new task using only a few datapoints and training iterations.



**Adapting to new sentence functions?**

[Credit: Finn and Levine, 2019]

# Proposed Approach

## Problem Formulation

**A Single Task:** response generation conditioned on a query-response sentence function pair $(d_X, d_Y)$

**Training Data**: K high-resource tasks: $D_{train}^k = \{(X_n^k, Y_n^k, d_X^k, d_Y^k), n = 1...N\}, k = 1...K$

**Testing Data (Target):** T tasks with infrequent sentence function: $D_{target}^t = \{(X_n^t, Y_n^t, d_X^t, d_Y^t), n = 1...N'\}, t = 1...T \ N' \ll N$

**Training:**

$$f_\theta : X^k \times (d_X^k, d_Y^k) \to Y^k, k = 1...K$$

Model    Query    Response    Task

**Testing:**

$$f_{\theta*} = \arg\max_\theta \log p(f_\theta | D_{target}^t, f_{\theta_0})$$

Adapted Model    Trained Model

Findings of ACL: EMNLP 2020, Fine-Grained Sentence Functions for Short-Text Conversation

# Proposed Approach

## Base Model: C-Seq2Seq

We use a conditional sequence-to-sequence learning model as our base model.

- Attentional Sequence-to-Sequence Model

- Learn an additional query-response sentence function embedding for each query-response type

- The sentence function embedding is used at every decoding step:

$$\mathbf{u}_t = \text{LSTM}(\mathbf{u}_{t-1}, [\mathbf{w}_t; \mathbf{s}_k])$$

Word Emb  Sentence Function Emb

# Proposed Approach

## MAML for C-Seq2Seq

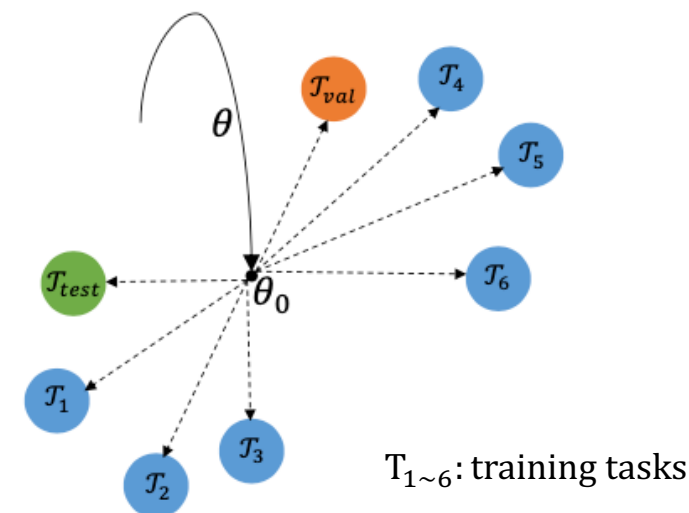Goal behind MAML: the conditions between task adaptation (fine-tuning) stage and training stage must match.



**Algorithm 2** Training of MAML

**Require:** $\mathcal{E}$: distribution over tasks $\{\mathcal{T}_1, ..., \mathcal{T}_K\}$
**Require:** $\alpha, \beta$: step size hyperparameters
1: Randomly initialize $\theta$
2: **while** not done **do**
3:     Sample a batch of tasks $\mathcal{T}_k \sim \mathcal{E}$
4:     **for** all $\mathcal{T}_k$ **do**
5:         Sample $D_{\mathcal{T}_k}, D'_{\mathcal{T}_k}$ from $\mathcal{T}_k$
6:         Evaluate $\nabla_{\theta_k} \mathcal{L}(f_{\theta_k})$ with respect to $D_{\mathcal{T}_k}$
7:         Update $\theta'_k = \theta_k - \alpha \nabla_{\theta_k} \mathcal{L}(f_{\theta_k})$
8:     **end for**
9:     Update $\theta \leftarrow \theta - \beta \sum_k \nabla_\theta \mathcal{L}(f_{\theta'_k})$ with respect to all $D'_{\mathcal{T}_k}$
10: **end while**

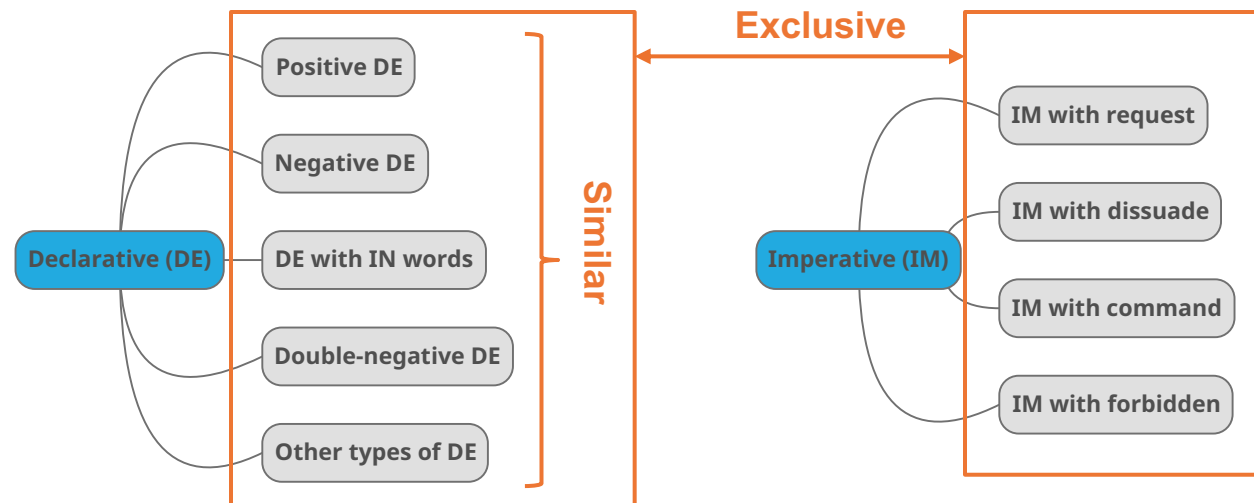$T_{1\sim6}$: training tasks

# Proposed Approach

## Exploring Structured Modeling

MAML assumes all tasks in training and adaptation stages distributed uniformly.

In conditioned response generation, some tasks may share some similarities while some are exclusive to each other.

# Proposed Approach

Exploring Structured Modeling: Structured Meta-Learning (SML)

Task Representation Learning: sentence function embeddings are used to interact with each other via a gated self-attention mechanism

Task-Specific Knowledge Adaptation: the self-attended representations of these sentence functions are used as parameter gates to tailor the transferable knowledge of the meta-learned prior parameters.

# Proposed Approach
## Task Representation Learning

1. Sentence functions (tasks) seen in training: $\mathbf{S} = [\mathbf{s}_1, \ldots, \mathbf{s}_K]$

2. Self-matching Operation:

$$\mathbf{a}_k = \text{softmax}(\mathbf{S}^\top \mathbf{s}_k), \quad \mathbf{m}_k = \mathbf{S}\mathbf{a}_k$$
$$\mathbf{f}_k = \tanh(\mathbf{W}_f[\mathbf{s}_k; \mathbf{m}_k]),$$

3. Gated summation for the final sentence function representation:

$$\mathbf{g}_k = \text{sigmoid}(\mathbf{W}_g[\mathbf{s}_k; \mathbf{m}_k])$$
$$\tilde{\mathbf{s}}_k = \mathbf{g}_k \odot \mathbf{f}_k + (1 - \mathbf{g}_k) \odot \mathbf{s}_k$$

4. $\tilde{\mathbf{s}}_k$ replaces $\mathbf{s}_k$ as input at each decoding time step
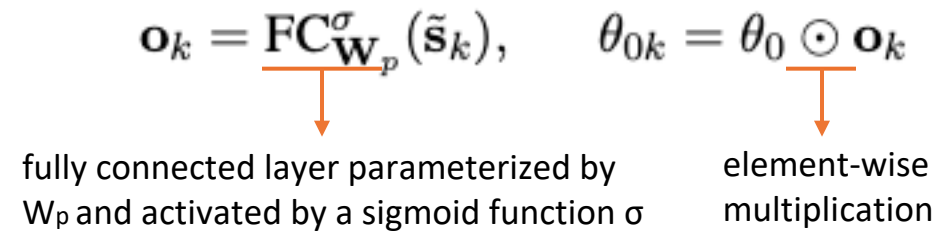
5. In the testing stage, new sentence functions can benefit from the learned sentence functions for fast adaptation.

Findings of ACL: EMNLP 2020, Fine-Grained Sentence Functions for Short-Text Conversation

# Proposed Approach

## Task-Specific Knowledge Adaptation

To adapt globally transferable knowledge $\theta_0$ to each sentence function, we design a parameter gate $\mathbf{o}_k$ for $\theta_0$:

$$\mathbf{o}_k = \mathrm{FC}^{\sigma}_{\mathbf{W}_p}(\tilde{\mathbf{s}}_k), \qquad \theta_{0k} = \theta_0 \odot \mathbf{o}_k$$

fully connected layer parameterized by $\mathrm{W_P}$ and activated by a sigmoid function σ

element-wise multiplication

Intuitively, sentence functions with similar representations will activate similar initial parameters while dissimilar sentence functions trigger different ones.
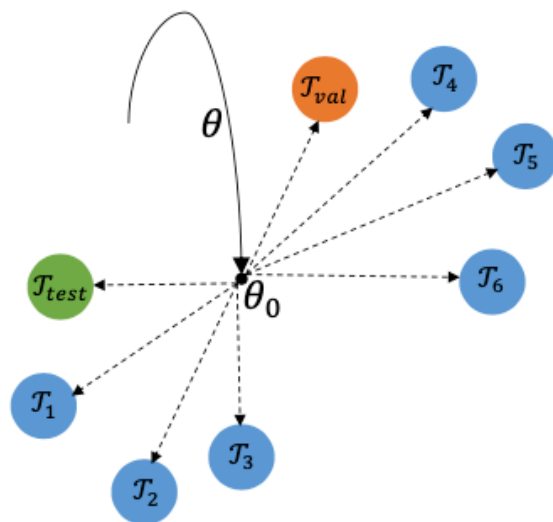
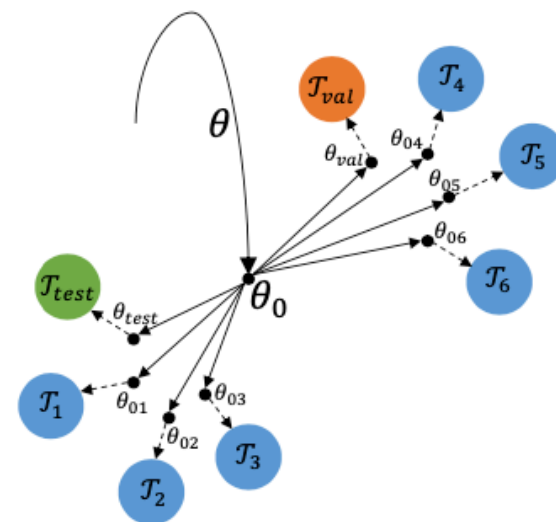# Proposed Approach

## MAML vs. SML

**Algorithm 3** Training of SML

**Require:** $\mathcal{E}$: distribution over tasks $\{\mathcal{T}_1, ..., \mathcal{T}_K\}$
**Require:** $\alpha, \beta$: step size hyperparameters
1: Randomly initialize $\theta$
2: **while** not done **do**
3:   Sample a batch of tasks $\mathcal{T}_k \sim \mathcal{E}$
4:   **for** all $\mathcal{T}_k$ **do**
5:     Sample $D_{\mathcal{T}_k}, D'_{\mathcal{T}_k}$ from $\mathcal{T}_k$
6:     Compute task representation $\tilde{\mathbf{s}}_k = \mathbf{g}_k \odot \mathbf{f}_k + (1 - \mathbf{g}_k) \odot \mathbf{s}_k$
7:     Compute $\mathbf{o}_k = \text{FC}^\sigma_{\mathbf{W}_p}(\tilde{\mathbf{s}}_k), \theta_{0k} = \theta_0 \odot \mathbf{o}_k$
8:     Evaluate $\nabla_{\theta_{0k}} \mathcal{L}(f_{\theta_{0k}})$ with respect to $D_{\mathcal{T}_k}$
9:     Update $\theta'_{0k} = \theta_{0k} - \alpha \nabla_{\theta_{0k}} \mathcal{L}(f_{\theta_{0k}})$
10:   **end for**
11:   Update $\theta \leftarrow \theta - \beta \sum_k \nabla_{\theta'_{0k}} \mathcal{L}(f_{\theta'_{0k}})$ with respect to all $D'_{\mathcal{T}_k}$
12: **end while**



(b) MAML

(c) SML

$T_{1\sim 6}$: training tasks

# Experiment

## Dataset

- STC-SeFun dataset (Bi et al., 2019)

- A large-scale Chinese short text conversation dataset with manually labeled sentence functions

- We select 9 high-resource tasks for meta-training, 4 tasks for meta-validation and 5 tasks for testing (adaptation).

| | Query SF | Response SF | # Samples | | |
|---|---|---|---|---|---|
| Meta Train | Positive DE | Positive DE | 27058 | | |
| | Wh-style IN | Positive DE | 12854 | | |
| | Positive DE | Negative DE | 5831 | | |
| | Negative DE | Positive DE | 4006 | | |
| | Positive DE | Wh-style IN | 3935 | | |
| | A-not-A IN | Positive DE | 3508 | | |
| | Wh-style IN | Negative DE | 3367 | | |
| | Yes-no IN | Positive DE | 3267 | | |
| | Negative DE | Negative DE | 2466 | | |
| Meta Val | Wh-style IN | DE w/ IN words | 271 | 100 | 500 |
| | Negative DE | Wh-style IN | 161 | 100 | 500 |
| | Positive DE | EX w/ interjections | 134 | 100 | 500 |
| | Positive DE | DE w/ IN words | 120 | 100 | 500 |
| Meta Test | Positive DE | Yes-no IN | 1314 | 100 | 500 |
| | Yes-no IN | Negative DE | 893 | 100 | 500 |
| | Positive DE | EX w/o tone words | 846 | 100 | 500 |
| | A-not-A IN | Negative DE | 684 | 100 | 500 |
| | Wh-style IN | Wh-style IN | 488 | 100 | 500 |

# Experiment

## Result

**Flue**: Fluency measures the grammatical correctness of responses (1-5)

**Rele**: Relevance measures whether the response is a relevant reply to the query (1-5)

**Info**: Informativeness evaluates whether the response provides any meaningful information with regard to the query (1-5)

**Accu**: Accuracy evaluates whether the response is coherent with the given response sentence function (0-1)
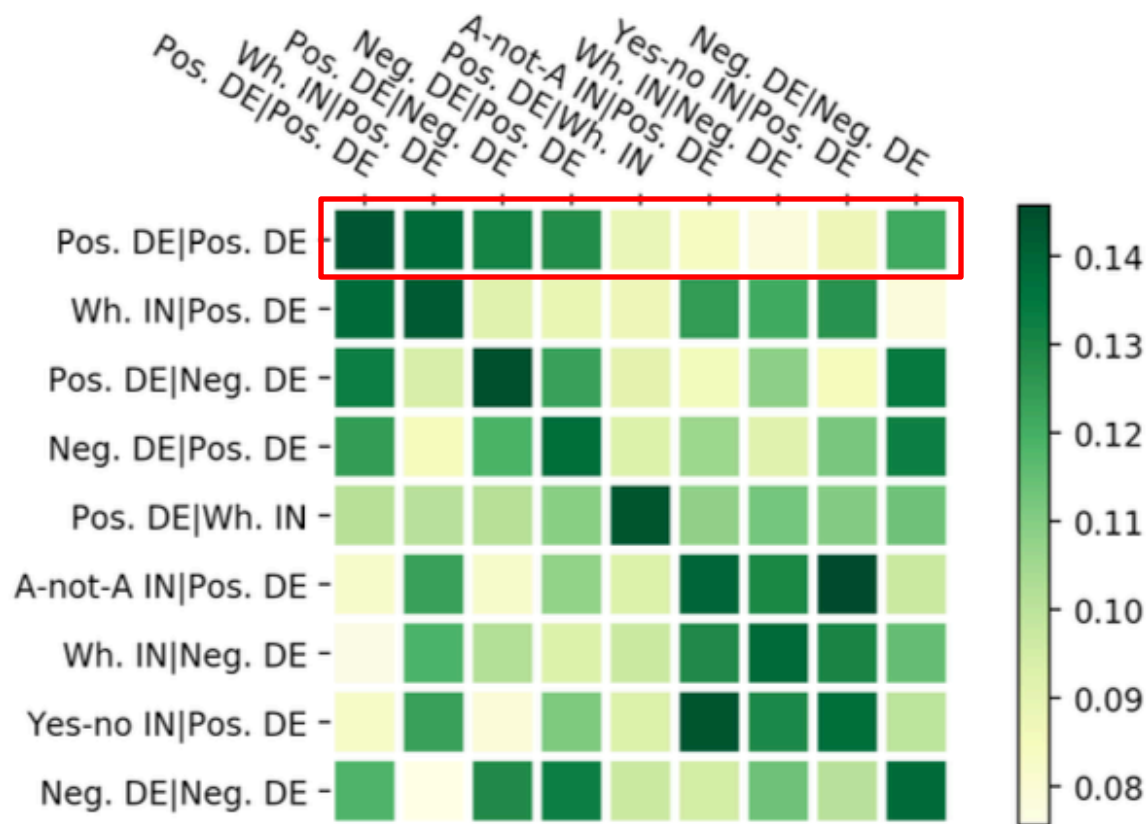
MTL: Multi-Task Learning

FT: Fine-tuning

SML: Structured Meta-Learning

| Query SF | Response SF | Model | Human Evaluation | | | |
|---|---|---|---|---|---|---|
| | | | Flue | Rele | Info | Accu |
| Positive DE | Yes-no IN | MTL | 59.40 | 50.40 | 36.53 | 0.33 |
| | | MTL+FT | 62.47 | 53.93 | 39.87 | 34.00 |
| | | MAML | 63.40 | 55.87 | 39.87 | 57.67 |
| | | SML | **64.27** | **56.00** | **40.13** | **69.00** |
| Yes-no IN | Negative DE | MTL | 60.47 | 57.07 | 49.87 | 6.00 |
| | | MTL+FT | 61.07 | 56.80 | 54.00 | 73.00 |
| | | MAML | 62.00 | **59.53** | 53.67 | **91.00** |
| | | SML | **64.93** | 57.80 | **55.93** | **91.00** |
| Positive DE | EX without tone words | MTL | 57.13 | 53.53 | 35.40 | 1.00 |
| | | MTL+FT | 56.40 | 53.67 | 36.67 | 39.00 |
| | | MAML | 65.33 | 56.20 | 39.27 | **71.00** |
| | | SML | **65.80** | **57.13** | **40.93** | 68.00 |
| A-not-A IN | Negative DE | MTL | 60.33 | 54.93 | 49.73 | 4.33 |
| | | MTL+FT | 62.13 | 55.13 | 51.47 | 53.67 |
| | | MAML | 62.60 | 55.20 | 51.27 | 89.33 |
| | | SML | **63.27** | **56.00** | **52.80** | **96.00** |
| Wh-style IN | Wh-style IN | MTL | 62.47 | 51.67 | 38.33 | 1.00 |
| | | MTL+FT | 63.60 | 52.60 | 39.13 | 22.33 |
| | | MAML | 64.07 | 53.13 | 43.33 | 85.00 |
| | | SML | **64.13** | **53.80** | **45.20** | **88.00** |

# Experiment

## Effect of Structure Modeling



Heatmap of the self-attention weight matrix. Each row shows the attention distribution for a given query-response sentence function pair (denoted in "Query|Response" format).

Findings of ACL: EMNLP 2020, Fine-Grained Sentence Functions for Short-Text Conversation

# Conclusion

- We apply model-agnostic meta-learning (MAML) for open domain dialogue generation on infrequent sentence functions.

- We further explore the structure across fine-grained sentence functions and such that the model can balance knowledge generalization and knowledge customization.

- Extensive experiments show that our structured meta-learning (SML) algorithm outperforms existing approaches under the low-resource setting.

# Thanks!