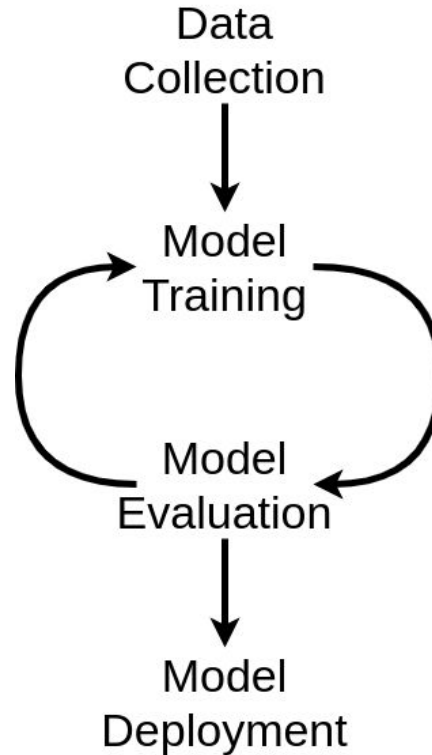


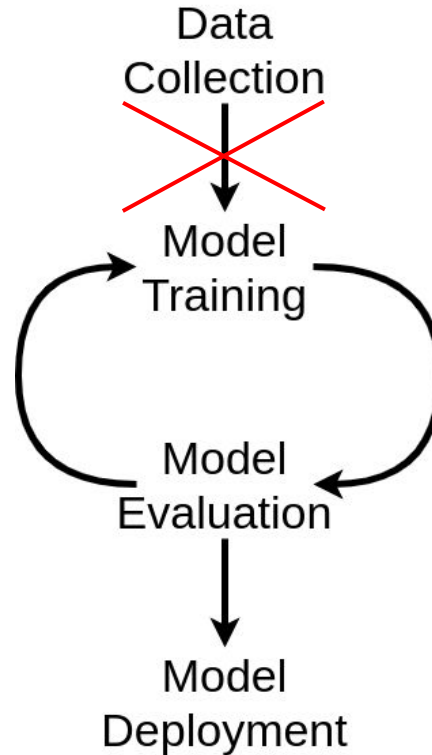
# Improving Data Efficiency for Natural Language Processing

Alon Albalak  
Ph.D. Proposal  
02/03/2023

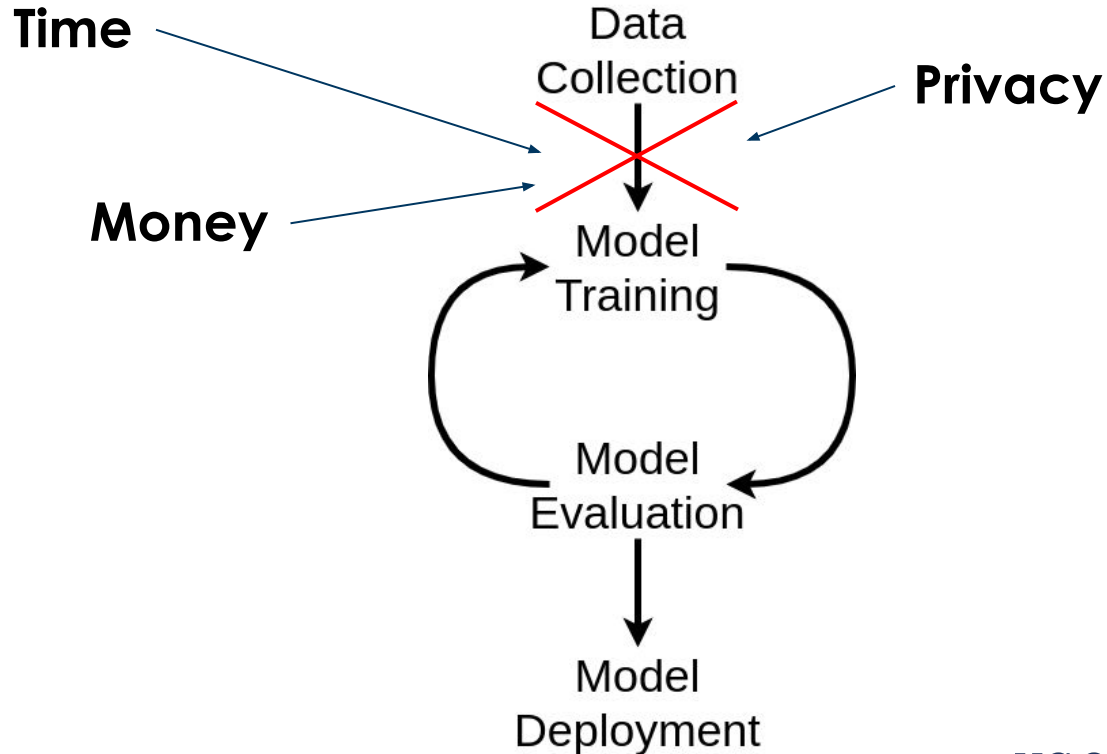
# Machine Learning Pipeline



# Machine Learning Pipeline

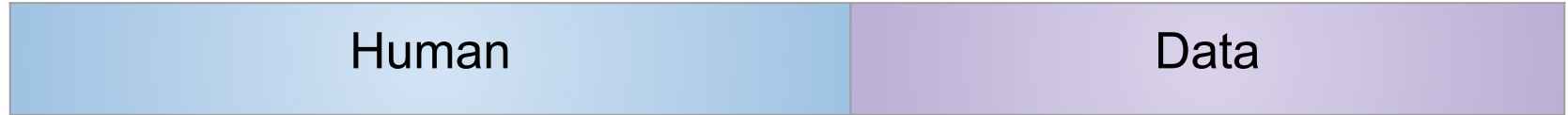


# Machine Learning Pipeline



# How can we utilize machine learning when DATA is a severe bottleneck?

# Transferring External Knowledge



# Transferring External Knowledge

Human	Data
<ul style="list-style-type: none"><li>● Model inductive bias</li><li>● Domain knowledge</li><li>● Relational information</li><li>● Data augmentation</li></ul>	

# Transferring External Knowledge

Human	Data
<ul style="list-style-type: none"><li>● Model inductive bias</li><li>● Domain knowledge</li><li>● Relational information</li><li>● Data augmentation</li></ul>	<ul style="list-style-type: none"><li>● Self-supervised domain data</li><li>● Labeled out-of-domain data</li><li>● Labeled related task data</li><li>● Data augmentation</li></ul>



# Self-supervised Training on Web Text

# Self-supervised Training on Web Text

GPT-3

w

The thesis proposal is a milestone on the way to receiving a Ph.D. The first milestone is the Major Area Exam (MAE), an oral presentation where the student presents their knowledge on their area of research. At the second milestone, the thesis proposal, the student prepares an oral presentation that describes the dissertation topic, summarizes the background literature, and presents a comprehensive research plan for the doctoral dissertation. The final exam is the defense of the student's dissertation, which consists of a public seminar and an evaluation by the student's doctoral committee on whether the student has successfully defended the dissertation.

Write a joke to be used at the students thesis proposal.



Why did the PhD student bring a ladder to their thesis proposal?



So they could reach the next level of academia!

# Self-supervised Training on Web Text

## GPT-3



Actually, can you write a knock-knock joke instead



Sure! Here's a PhD-themed knock-knock joke:



Knock, knock.

Who's there?

Thesis.

Thesis who?

Thesis what I've been working on for the past 5 years!

# My Research on Data Efficiency

Human	Data
<ul style="list-style-type: none"><li>● Albalak et al. 2022 (preprint)</li><li>● Augustine et al. 2022 (preprint)</li><li>● Dickens et al. 2022 (TPM workshop)</li><li>● Pryor et al. 2022 (preprint)</li></ul>	<ul style="list-style-type: none"><li>● Albalak et al. 2022 (EMNLP 2022)</li><li>● Albalak et al. 2022 (ENSLP workshop 2022)</li><li>● Li et al. 2022 (Alexa Prize 2022)</li><li>● Albalak et al. 2022 (ConvAI workshop 2022)</li></ul>

# My Research on Data Efficiency

Human	Data
<ul style="list-style-type: none"><li>● Albalak et al. 2022 (preprint)</li><li>● Augustine et al. 2022 (preprint)</li><li>● Dickens et al. 2022 (TPM workshop)</li><li>● Pryor et al. 2022 (preprint)</li></ul>	<ul style="list-style-type: none"><li>● Albalak et al. 2022 (EMNLP 2022)</li><li>● Albalak et al. 2022 (ENSLP workshop 2022)</li><li>● Li et al. 2022 (Alexa Prize 2022)</li><li>● Albalak et al. 2022 (ConvAI workshop 2022)</li></ul>

# Zero-shot Transfer Methods

**Problem:** No target data samples

**Idea:** Convert many tasks into text-to-text format

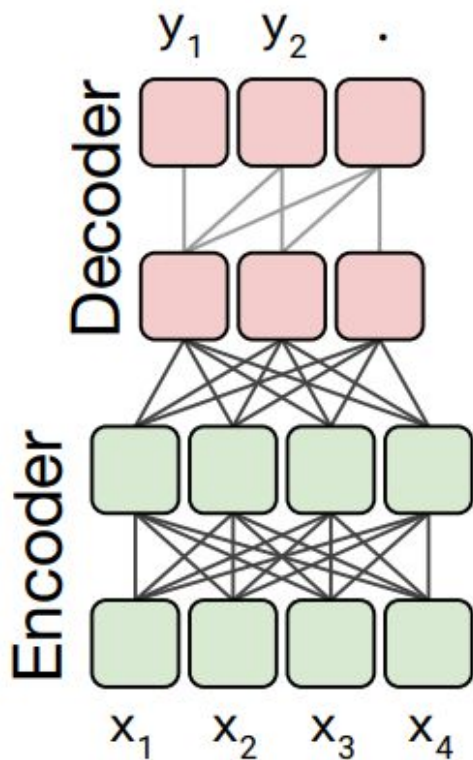
**Allows:** 1 model can perform multiple tasks

**An Exploration of Methods for Zero-shot Transfer in Small Language Models.**

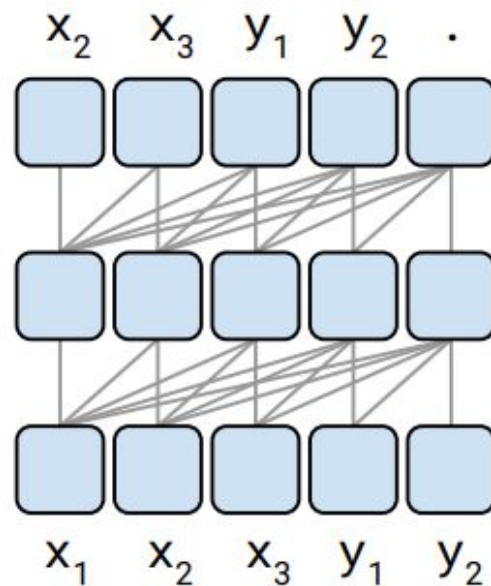
Alon Albalak, Akshat Shrivastava, Chinnadhurai Sankar, Adithya Sagar, Mike Ross.

*Efficient Natural Language and Speech Processing, 2022.*

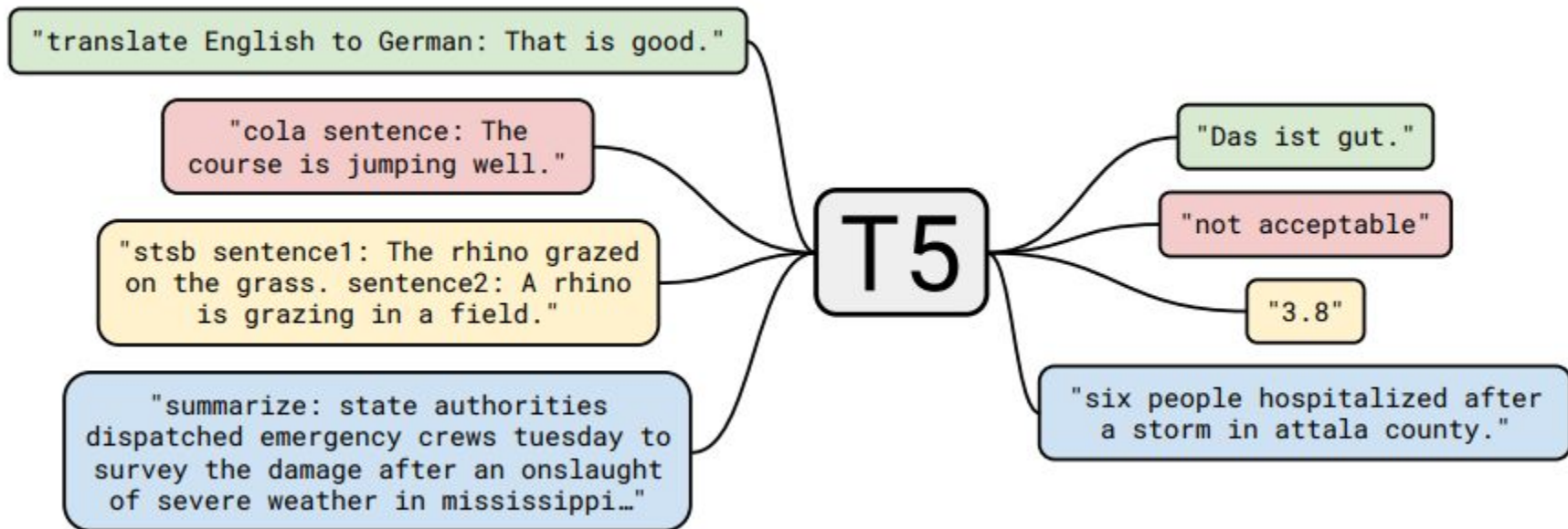
# Generative Models



## Language model

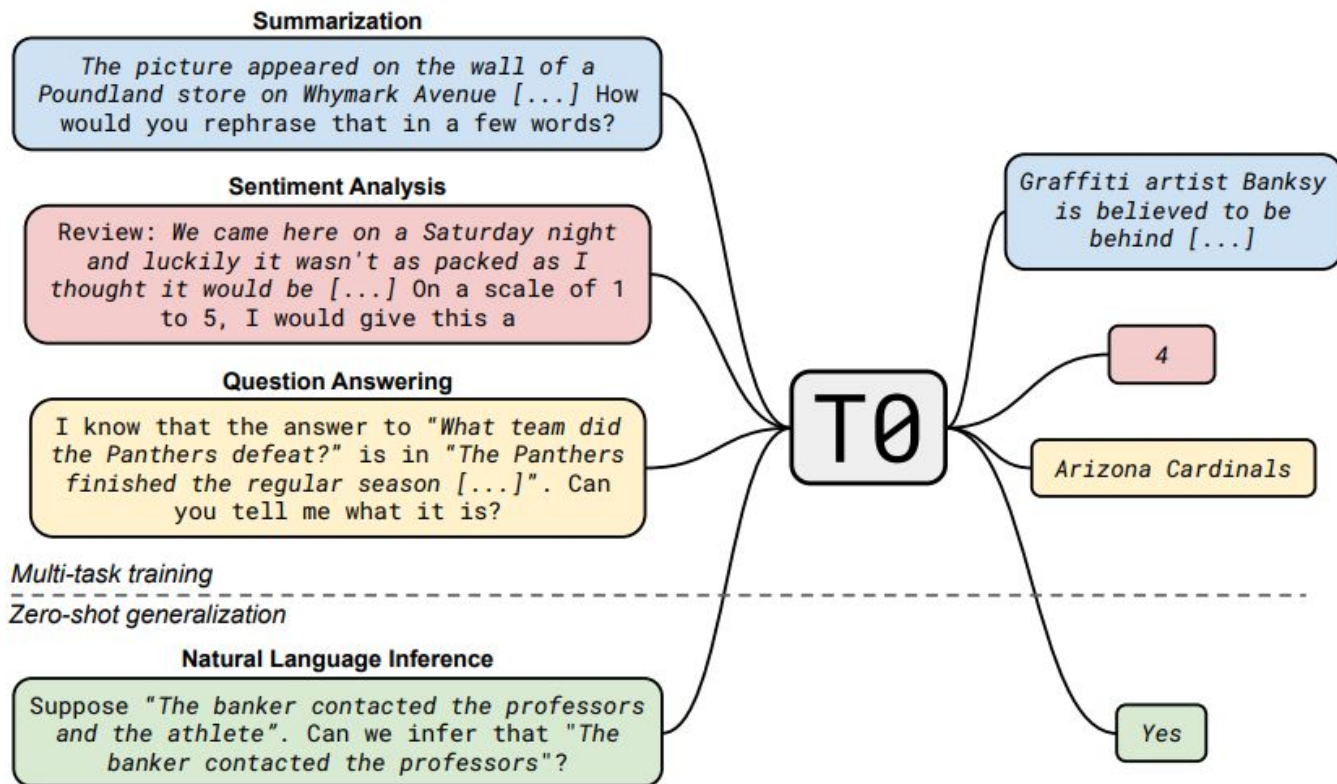


# Multi-tasking with Text-to-text Transfer





# Multi-tasking for Zero-shot Transfer



# Instruction Tuning

## Relation classification

**Instruction:** You will be given some conversation text and you need to find the relation in the conversation between specified people or speakers.

**Input:** [CONTEXT] Speaker 1: You know Phoebe, when I was little... [ENDOFTURN]  
Speaker 2: Oh, I love family ... [ENDOFDIALOGUE] Possible relations are:  
[OPTIONS] 0: students, 1: visited place, 3: schools attended, 4: siblings, ...  
[QUESTION] Choose the most possible relation between Speaker 1 and Ursula

↓  
**Output**

4

**Generalization methods have been well studied in large language models.**

**How do they interact within smaller language models?**

# Experiment Design

## Data

- 46 Tasks from InstructDial<sup>1</sup>
- Tasks have between 3 and 10 instructions
- 3 splits of train/test tasks
  - 40 train tasks
  - 6 test tasks
- Tasks are divided into classification and generation

## Models

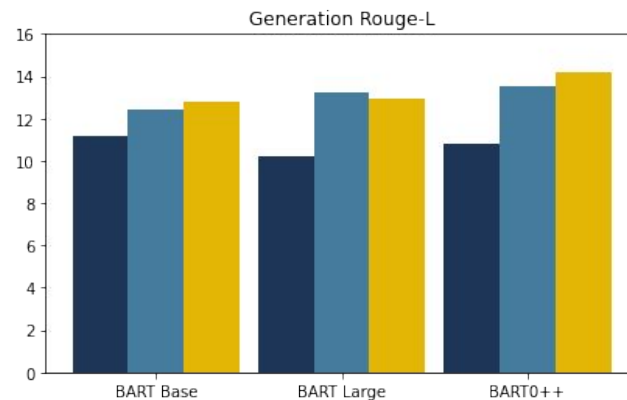
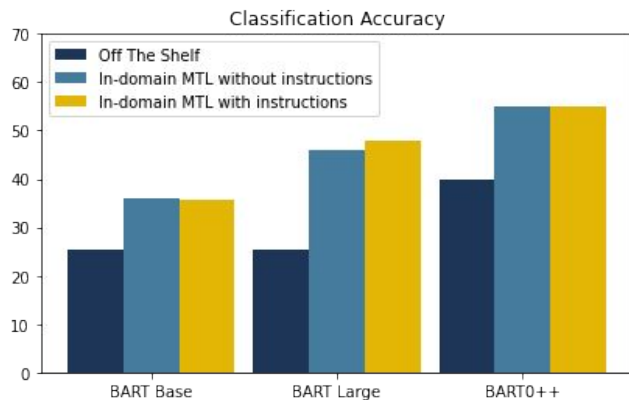
3 variants of BART

- BART-Base
- BART-Large
- BART0++

<sup>1</sup>Gupta et al. InstructDial: Improving Zero and Few-shot Generalization in Dialogue through Instruction Tuning. 2022

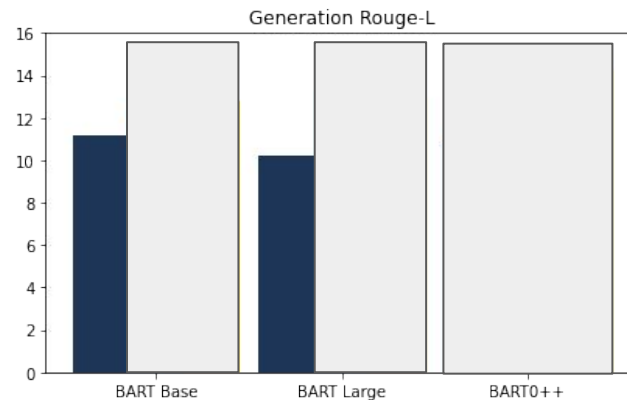
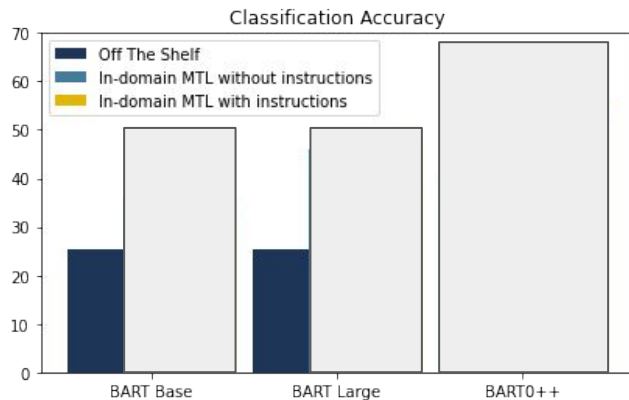
# Zero-Shot Results

- Model Size
- Multi-Task Learning (MTL)
- In-Domain (Dialogue) MTL
- Instruction Tuning



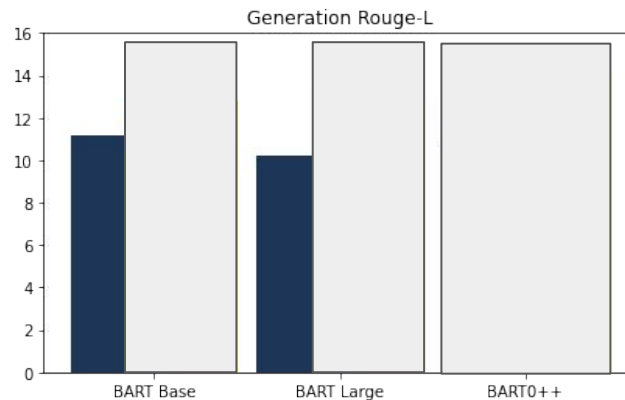
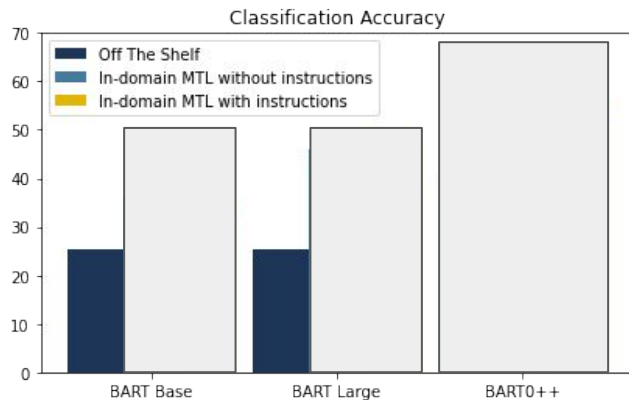
# Effects of Model Size

- Nearly identical performance on classification
- Slightly better BART-Base on generation



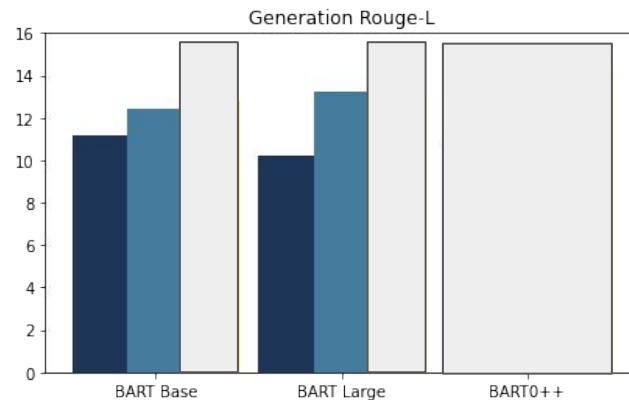
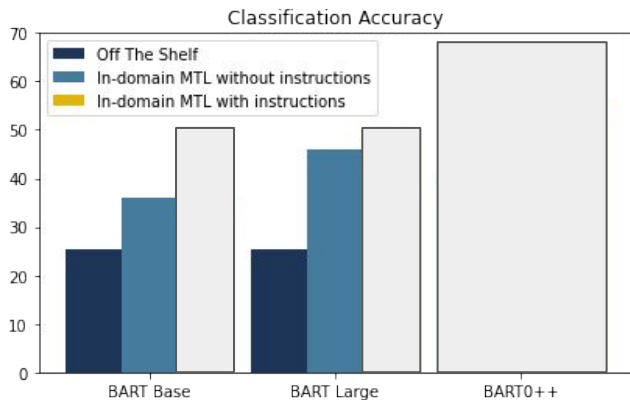
# Effects of Model Size

- Nearly identical performance on classification
- Slightly better BART-Base on generation
- Takeaway: With same data, larger model doesn't necessarily improve performance



# Effects of Model Size + In-domain MTL

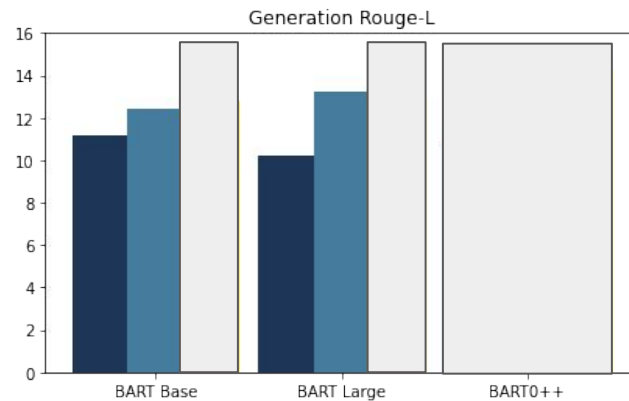
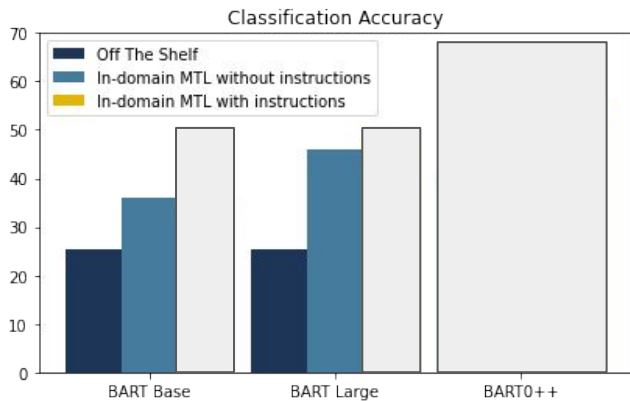
- BART-Base improves by 6.5 average
- BART-Large improves by 13.3 average





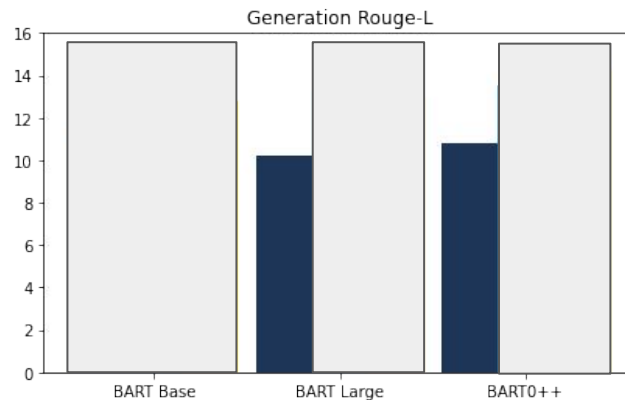
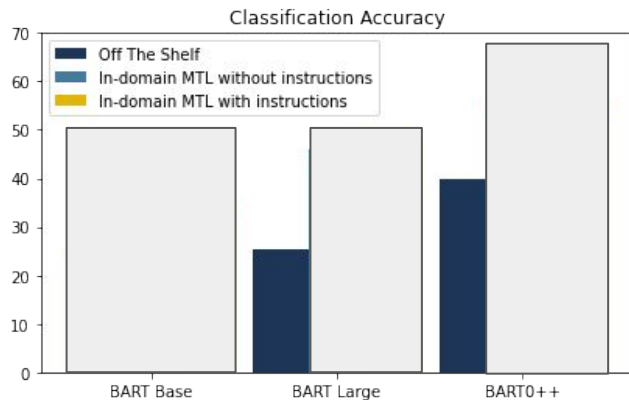
# Effects of Model Size + In-domain MTL

- BART-Base improves by 6.5 average
- BART-Large improves by 13.3 average
- Takeaway: Increasing model size AND training on in-domain data has better potential



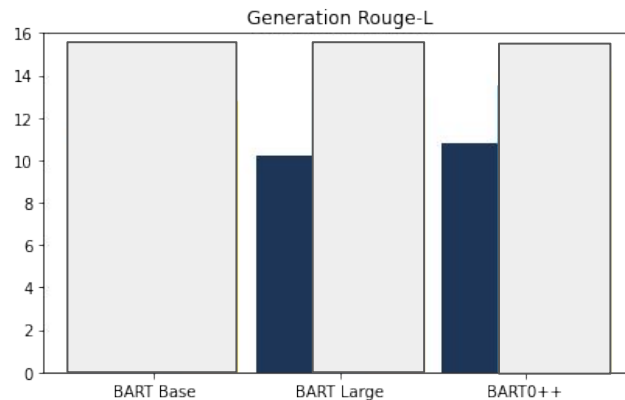
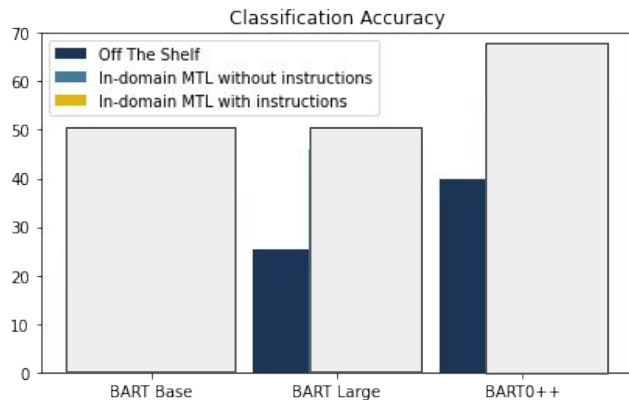
# Effect of General Purpose Multi-Task Learning

- 14.5 point (57.1% relative) improvement on classification
- 0.6 Rouge-L (5% relative) improvement on generation
- BART0++ MTL tasks are mainly classification



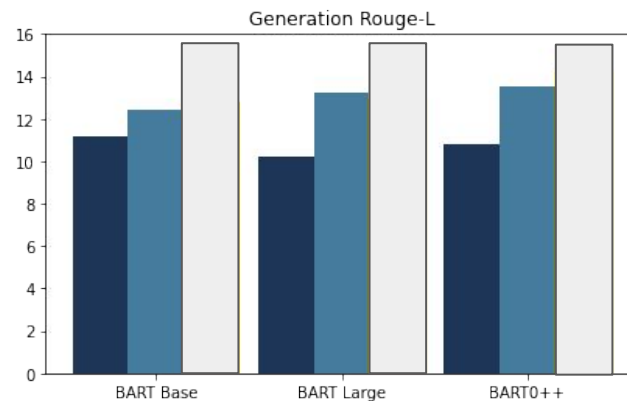
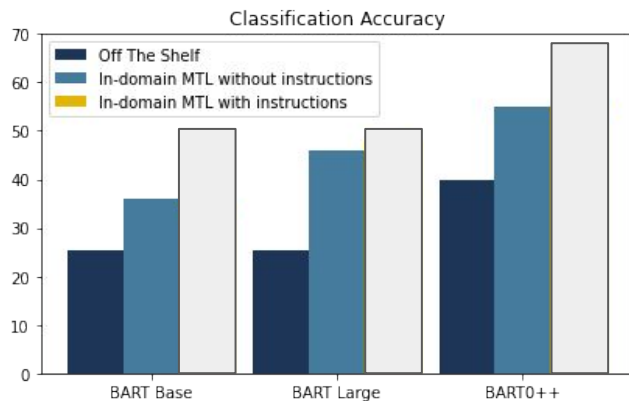
# Effect of General Purpose Multi-Task Learning

- 14.5 point (57.1% relative) improvement on classification
  - 0.6 Rouge-L (5% relative) improvement on generation
  - BART0++ MTL tasks are mainly classification
- Takeaway: General purpose MTL is incredibly beneficial when test tasks are in same distribution as MTL tasks



# Effects of In-Domain Multi-Task Learning

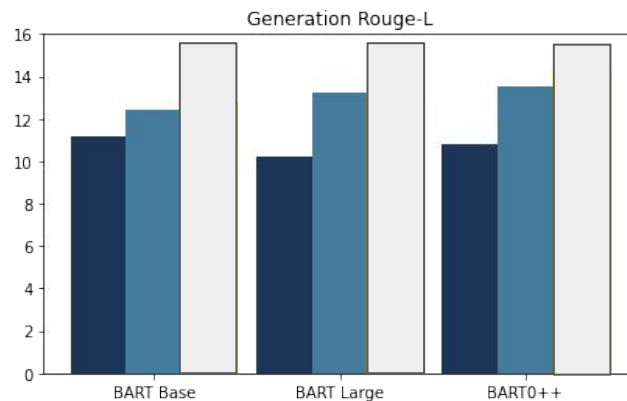
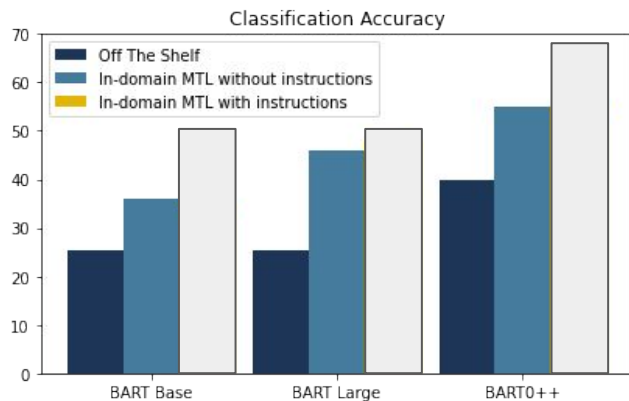
Relative Improvement	BART-Base	BART-Large	BART0++
Classification	41.8%	<b>80%</b>	37.7%
Generation	11.5%	<b>29.3%</b>	25.3%



# Effects of In-Domain Multi-Task Learning

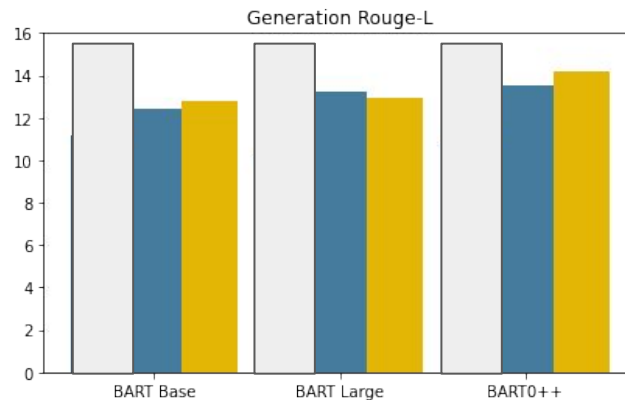
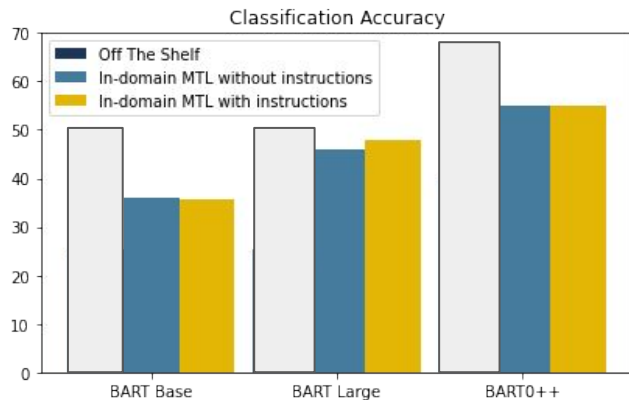
Relative Improvement	BART-Base	BART-Large	BART0++
Classification	41.8%	<b>80%</b>	37.7%
Generation	11.5%	<b>29.3%</b>	25.3%

➤ Takeaway: In-domain MTL gives largest portion of generalization improvement



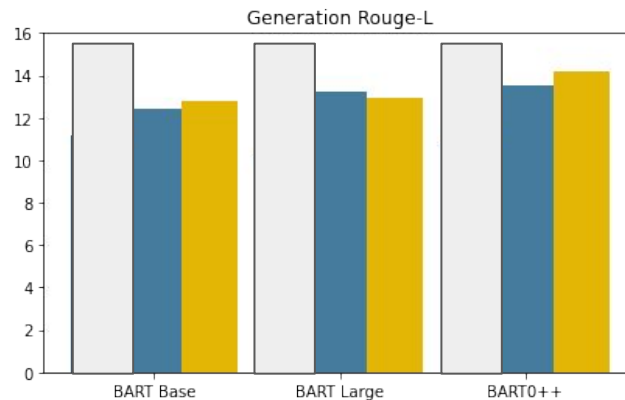
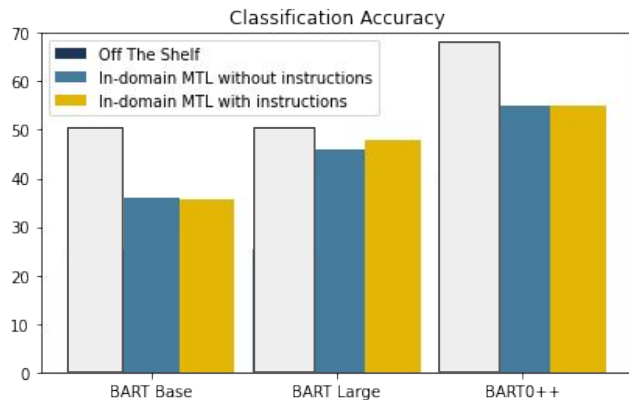
# Effects of Instruction Tuning

- Much smaller effect than previous variables
- 2% average improvement across models
- This finding runs counter to previous studies which found ~10% decrease in performance for models w/ <8B parameters



# Effects of Instruction Tuning

- Much smaller effect than previous variables
  - 2% average improvement across models
  - This finding runs counter to previous studies which found ~10% decrease in performance for models w/ <8B parameters
- Takeaway: Instructions are beneficial, but gains are diminishing



**Small models can benefit from multi-task training for zero-shot generalization**

**If we have a specific target task in mind, can we identify which tasks will transfer well?**



# Previous Studies on Task Transfer

Name	Train	Dev	task	metrics	genre/source
CommonsenseQA	9,741	1,221	question answering	acc.	ConceptNet
SciTail	23,596	1,304	natural language inference	acc.	science exams
Cosmos QA	25,588	3,000	question answering	acc.	blogs
SocialQA	33,410	1,954	question answering	acc.	crowdsourcing
CCG	38,015	5,484	tagging	acc.	Wall Street Journal
HellaSwag	39,905	10,042	sentence completion	acc.	video captions & Wikihow
QA-SRL	44,837	7,895	question answering	F1/EM	Wikipedia
SST-2	67,349	872	sentiment classification	acc.	movie reviews
QAMR	73,561	27,535	question answering	F1/EM	Wikipedia
QQP	363,846	40,430	paraphrase detection	acc./F1	Quora questions
MNLI	392,702	20,000	natural language inference	acc.	fiction, letters, telephone speech
CB	250	57	natural language inference	acc./F1	Wall Street Journal, fiction, dialogue
COPA	400	100	question answering	acc.	blogs, photography encyclopedia
WSC	554	104	coreference resolution	acc.	hand-crafted
RTE	2,490	278	natural language inference	acc.	news, Wikipedia
MultiRC	5,100	953	question answering	F1 <sub>α</sub> /EM	crowd-sourced
WiC	5,428	638	word sense disambiguation	acc.	WordNet, VerbNet, Wiktionary
BoolQ	9,427	3,270	question answering	acc.	Google queries, Wikipedia
CommonsenseQA	9,741	1,221	question answering	acc.	ConceptNet
Cosmos QA	25,588	3,000	question answering	acc.	blogs
ReCoRD	100,730	10,000	question answering	F1/EM	news (CNN, Daily Mail)

## Classification

### Sentiment Analysis

Amazon\_Polarity (McAuley et al. 2013)  
 IMDB (Maas et al. 2011)  
 Poem\_Sentiment (Sheng et al. 2020) ...

### Paraphrase Identification

Quora Question Paraphrases (Quora)  
 MRPC (Dolan et al. 2005)  
 PAWS (Zhang et al. 2019) ...

### Natural Language Inference

MNLI (Williams et al. 2018)  
 QNLI (Rajpurkar et al. 2016)  
 SciTail (Knot et al. 2018) ...

Others (topic, hate speech, ...)

## Question Answering

### Reading Comprehension

SQuAD (Rajpurkar et al. 2016)  
 QuoRef (Dasigi et al. 2019)  
 TweetQA (Xiong et al. 2019) ...

### Multiple-Choice QA

CommonsenseQA (Talmor et al. 2019)  
 OpenbookQA (Mihaylov et al. 2018)  
 AI2\_ARC (Clark et al. 2018) ...

### Closed-book QA

WebQuestions (Berant et al. 2013)  
 FreebaseQA (Jiang et al. 2019)  
 KILT-NQ (Kwiatkowski et al. 2019) ...

Others (yes/no, long-form QA)

## Conditional Generation

### Summarization

Gigaword (Napoles et al. 2012)  
 XSum (Narayan et al. 2018) ...

### Dialogue

Empathetic Dialog (Rashkin et al. 2019)  
 KILT-Wow (Dinan et al. 2019) ...

Others (text2SQL, table2text ...)

## Others

### Regression

Mocha (Chen et al. 2020)  
 Yelp Review Full (Yelp Open Dataset) ...

### Others

Acronym Identification  
 Sign Language Translation  
 Autoregressive Entity Linking  
 Motion Recognition  
 Pronoun Resolution ...

# Intra-Dataset Task Transfer With FETA

## *Intra-Dataset Task Transfer:*

Transferring knowledge from a source task to a target task, where both source and target are in the same distribution (domain)

### **FETA: A Benchmark for Few-Sample Task Transfer in Open-Domain Dialogue**

Alon Albalak, Yi-Lin Tuan, Pegah Jandaghi, Connor Pryor, Luke Yoffe, Deepak Ramachandran, Lise Getoor, Jay Pujara, William Yang Wang.

*EMNLP 2022*

# Intra-Dataset Task Transfer With FETA-Friends

## Dialogue

Chandler (C), Rachel (R), Gunther (G)

- (1) C: (Reading a comic) Eh, I don't know.
- (2) R: What?
- (3) C: Well, as old as he is in dog years, do you think Snoopy should still be allowed to fly this thing?
- (4) G: Rachel?
- (5) R: Yeah.
- (6) G: Do you remember when you first came here, how you spent two weeks getting trained by another waitress?
- (7) R: Oh sure. Do you need me to train somebody?
- (8) G: (laughs) Good one. Actually, Terry wants you to take the training again.
- (9) R: (To Chandler) Eh, do you believe that?
- (10) C: (Thinks about it) Yeah.

## Tasks

### MELD Emotion Recognition

Answer: utt(6) = Neutral

### Question Answering

Question: How long did Rachel train for?  
Answer: two weeks

### Emory Emotion Recognition

Answer: utt(6) = Powerful

### Personality Detection

Subject: Gunther  
Answer:  
Agreeable = Yes  
Conscientious = No  
Extroverted = Yes  
Open = Yes  
Neurotic = Yes

### Reading Comprehension

Statement: Gunther interrupts Rachel talking to [??] while on the job and says Terry needs the new-waitress training to be taken again.  
Question: Out of Chandler, Rachel, Gunther, and Terry who is [??] ?  
Answer: Contradiction

### Relation Extraction

(Head, Tail) -> Relation:  
(Rachel, Waitress) -> hasTitle  
(Rachel, Terry) -> hasBoss

### Character Identification

Question: In utt (3), who does 'he' refer to?  
Answer: "Snoopy"

## FETA: A Benchmark for Few-Sample Task Transfer in Open-Domain Dialogue

Alon Albalak, Yi-Lin Tuan, Pegah Jandaghi, Connor Pryor, Luke Yoffe, Deepak Ramachandran, Lise Getoor, Jay Pujara, William Yang Wang.  
EMNLP 2022

# Intra-Dataset Task Transfer With FETA-DailyDialog

## Dialogue

- (1) **A:** Happy anniversary, sweetheart!  
(2) **B:** Yes. to our first anniversary and many more to come. Cheers!  
(3) **A:** I'll drink to that. Thanks for making this a night worth remembering.  
(4) **B:** Well, it's a special day. They say if you survive the first year, the rest is smooth sailing.  
(5) **A:** That's good to know. Oh, listen! The band's playing our song.  
(6) **B:** I requested it. What do you say? Do you have your dancing shoes on?  
(7) **A:** Always.

## Tasks

**Topic Classification**  
Answer: Relationship

**Causal Emotion Span Extraction**  
Question: In utt(2), what causes the emotion "happiness"?  
Answer: "first anniversary"

**Causal Emotion Entailment**  
Premise: utt(1) - utt(5)  
Hypothesis: utt(3) causes "happiness" in utt(5)  
Answer: Contradiction

**Dialogue-Level NLI**  
Premise: utt(1) - utt(7)  
Hypothesis: "dancing" happens simultaneously with "song"  
Answer: Entailment

**Dialog Act Classification**  
Answer: utt(1) = Inform

**Dialogue Reasoning Multiple Choice**  
Question: From utt(3), when does "drink to that" happen?  
Options:  
(a) Survive the first year  
(b) Anniversary  
(c) Night  
(d) Playing our song  
Answer: Option (b)

**Commonsense Relation Extraction**  
Question: What is the relationship between "special day" and "night worth remembering"?  
Answer: "special day" causes "night worth remembering"

**Emotion Recognition**  
Answer: utt(1) = Happiness

**Dialogue Reasoning Span Extraction**  
Context: utt(1) - utt(6)  
Response Options:  
(a) I think I should skip this.  
(b) That's good to know, listening to music makes your day happy.  
(c) Yes, let's go to the floor.  
Answer: Option (c)

**Adversarial Response Selection**  
Question: In utt(2), what is "anniversary"?  
Answer: "special day"

## FETA: A Benchmark for Few-Sample Task Transfer in Open-Domain Dialogue

Alon Albalak, Yi-Lin Tuan, Pegah Jandaghi, Connor Pryor, Luke Yoffe, Deepak Ramachandran, Lise Getoor, Jay Pujara, William Yang Wang.  
EMNLP 2022

# FETA Learning Setting

- Pairwise task transfer within dataset
- 2 sets with 10 and 7 tasks = 132 source-target pairs
- For each experiment:
  - Source task uses full data
  - Target task uses 10% of data

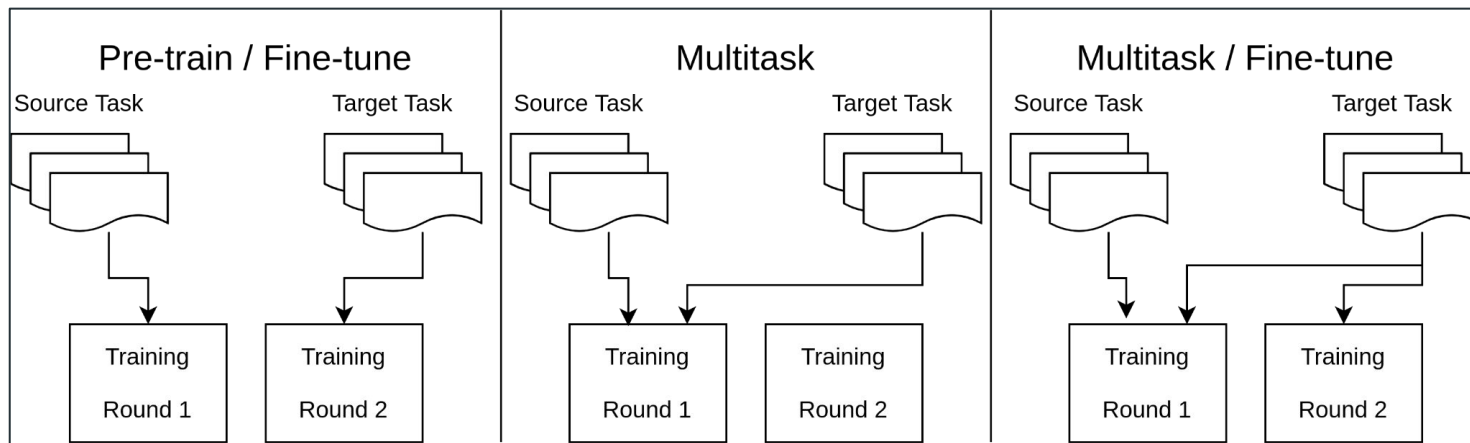
# Experiments

## 3 Model Architectures

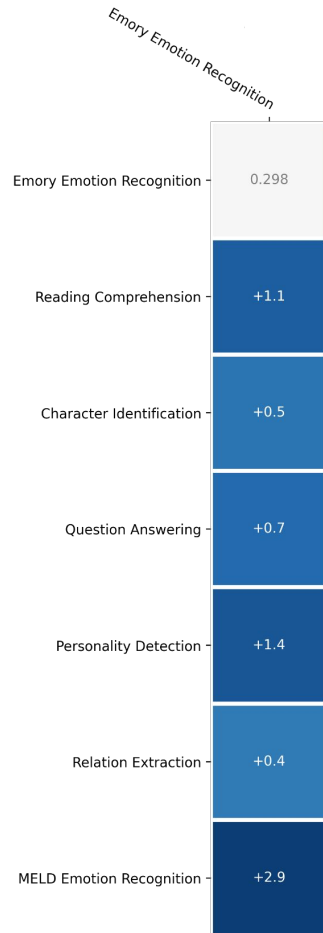
- Encoder - BERT
- Decoder - GPT
- Encoder-Decoder - T5

## 3 Transfer Algorithms

- Pre-train/Fine-tune
- Multitask
- Multitask/Fine-tune



# Sample Experiment

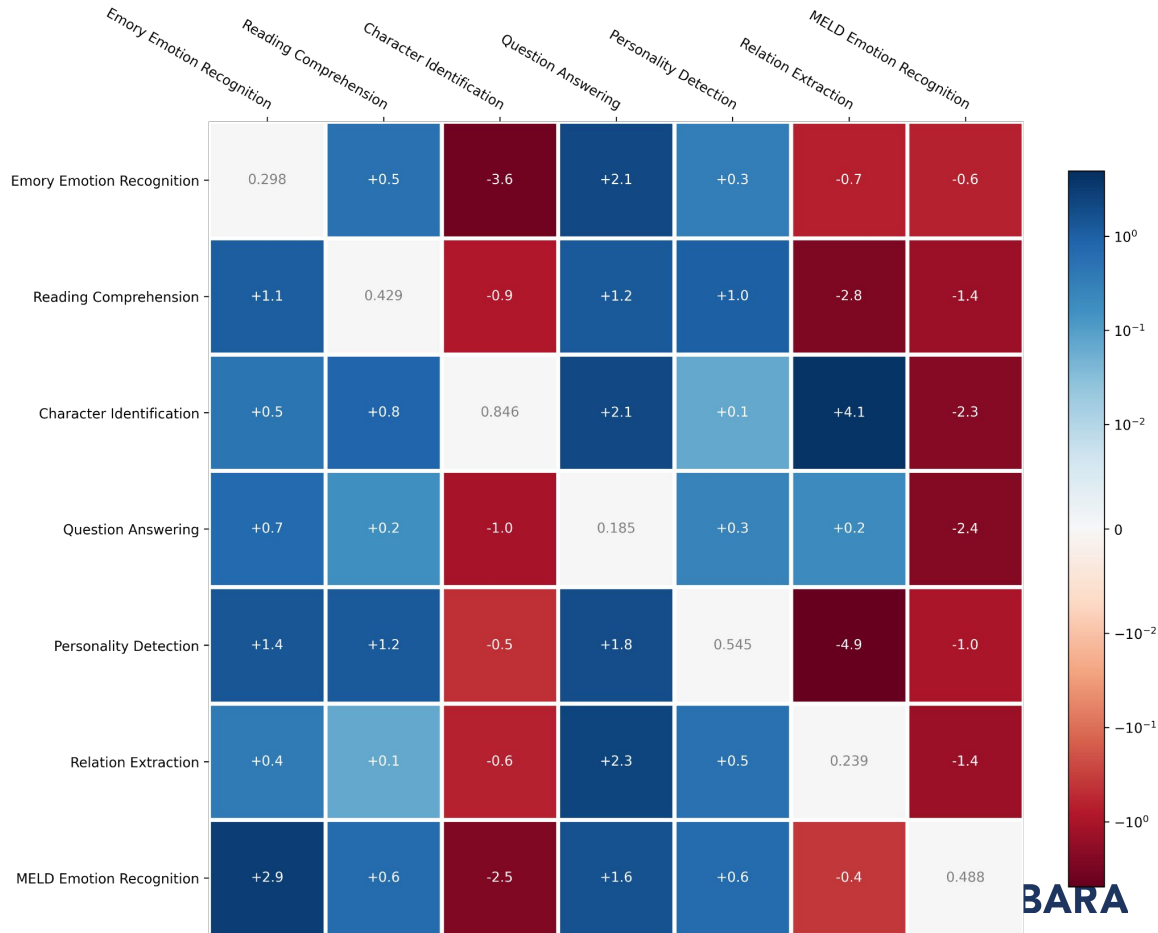


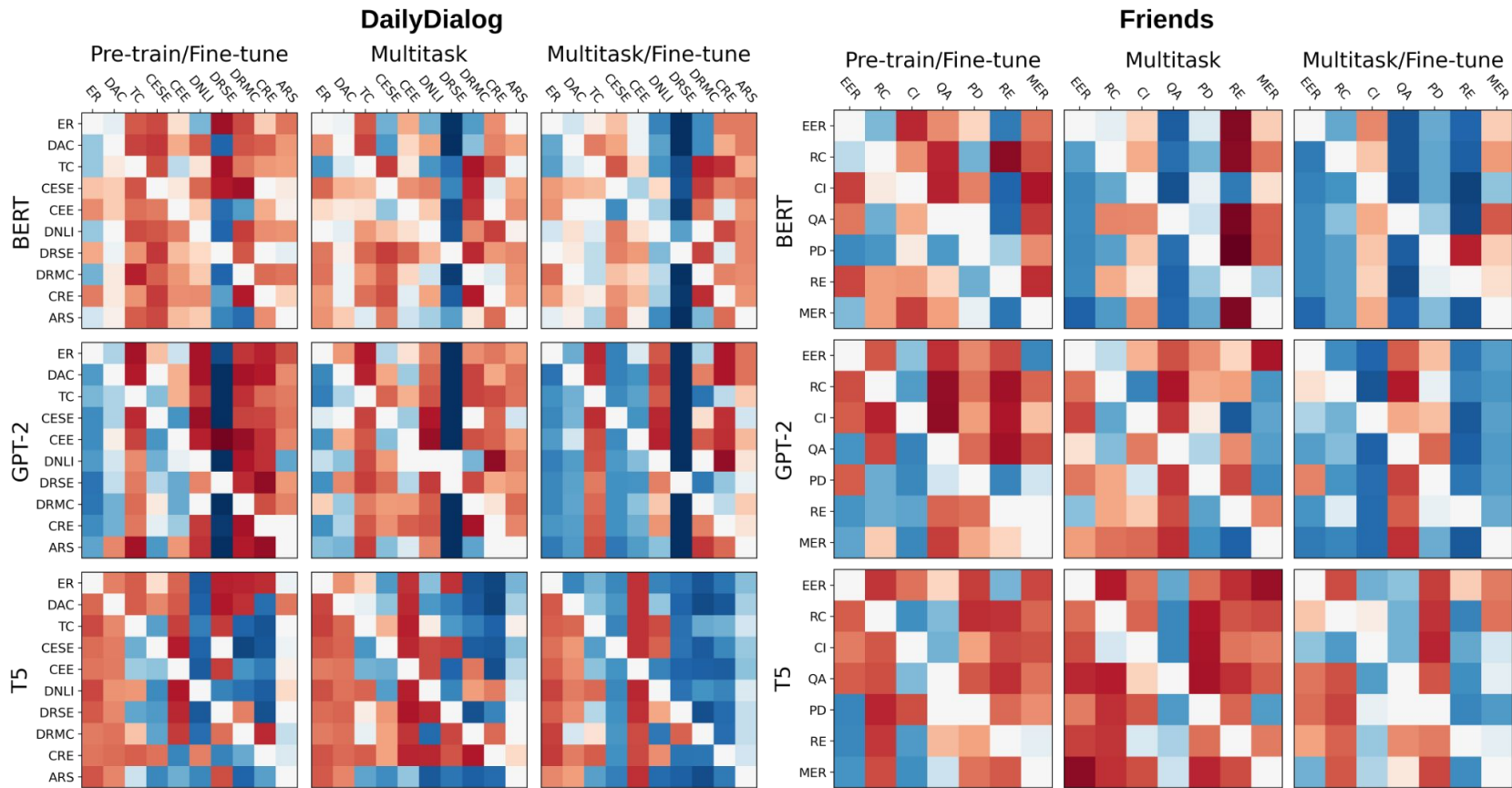
# Sample Experiment





# Sample Experiment





# FETA Takeaways

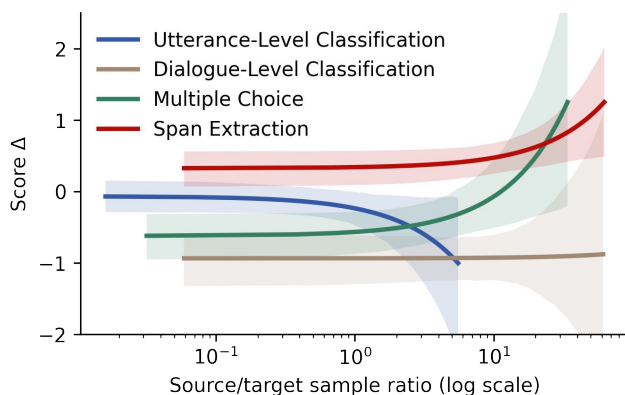
1. Finding the best source task can make a big difference
2. Label-space complexity affects transfer
3. Mitigating negative transfer with multitasking
4. Adding multiple sources can negatively impact transfer

# Takeaway - Gap between average and best source

- We find the difference between using the best source task vs. average of all source tasks to be  $\sim 1.6$  points, with the largest gap being 3.5 points
- **This strongly motivates the need for further understanding which source tasks will transfer best to specific target tasks**

# Takeaway - Effect of label-space complexity

- We find that span extraction target tasks have positive transfer from all source task types
- Multiple choice target tasks also see positive transfer, but only when the ratio of source-to-target samples is large ( $>10$ -to-1 as shown in Fig. 2 below)
- Both classification tasks see increasingly negative transfer with increasing number of source task samples
- **Overall, as the label-space of a target task becomes more complex, the task benefits more from transfer**



## Takeaway - Mitigate negative transfer with Multitask/Fine-tune

- We find that T5's best individual scores are with Pre-train/Fine-tune, but the best average scores are with Multitask/Fine-tune
- In fact, for GPT-2 on FETA-Friends, using the worst source task will still lead to a 0.74% improvement over the baseline
- **We find that across algorithms, Multitask/Fine-tune achieves the best worst-case performance**

## Takeaway - Adding source tasks can negatively impact transfer

- Small scale experiment:
  - 4 target tasks with best source tasks
  - Train models using top-3 source tasks
- Results:
  - GPT-2 improves most (8/12 settings)
  - BERT improves on 5/12 settings
  - T5 only improves on 4/12 settings
- Takeaway - **Naively adding source tasks can actually hurt performance**

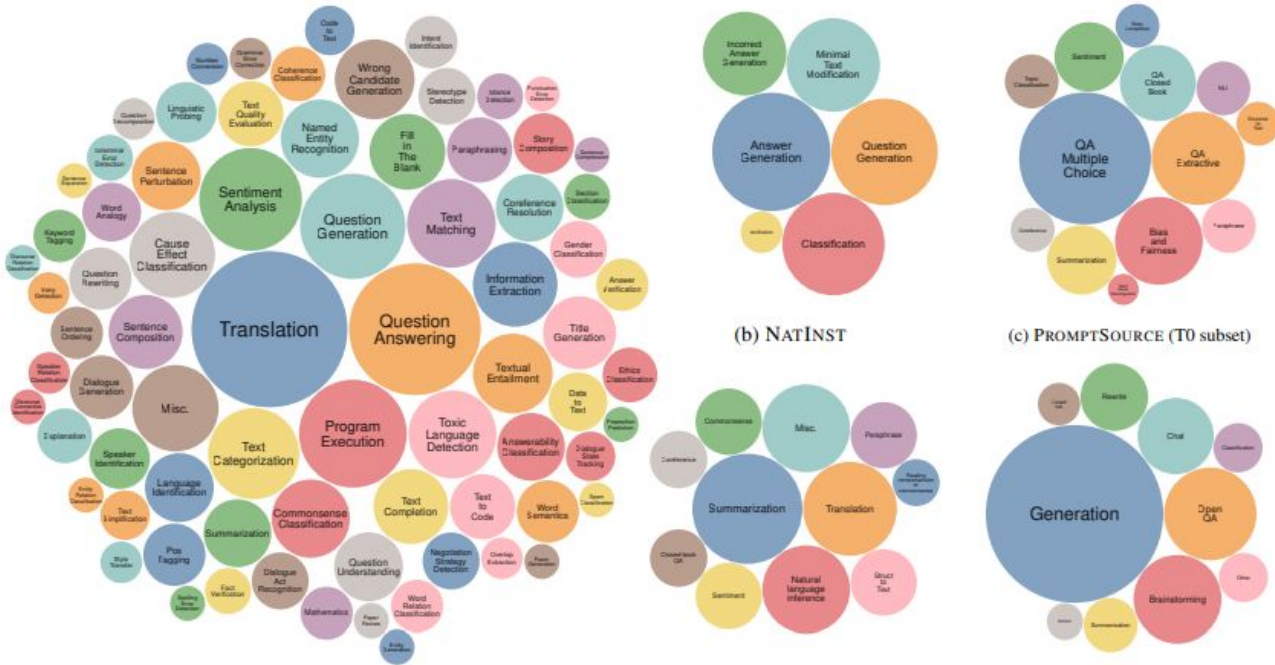
**Not all source tasks are equal and adding more source tasks doesn't always improve performance**

**Moving forward, how can we mitigate negative transfer?**

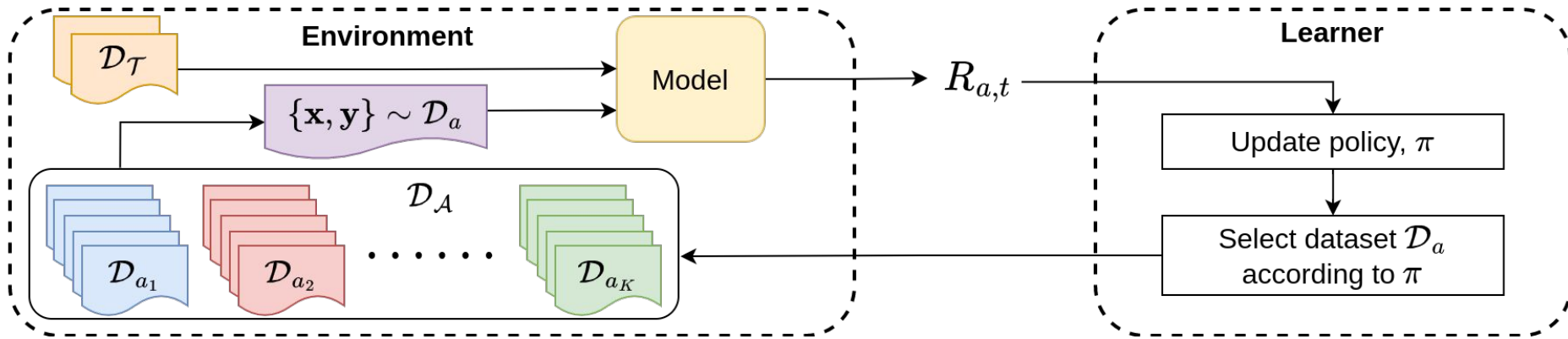


# Current work

## Few-shot Learning with Auxiliary Datasets - FLAD



# Few-shot Learning with Multi-armed Bandits



- Explore (auxiliary dataset) arms to find which gives best reward
- Exploit knowledge of previous rewards
- Use gradient alignment as reward

# Future Work

- Data pruning/selection
  - How can we determine which samples to best use, efficiently
- Demonstration selection for in-context learning (ICL)
  - How can we best select demonstrations for ICL based on a given task instance

# Questions

 @AlbalakAlon

 [alon\\_albalak@ucsb.edu](mailto:alon_albalak@ucsb.edu)

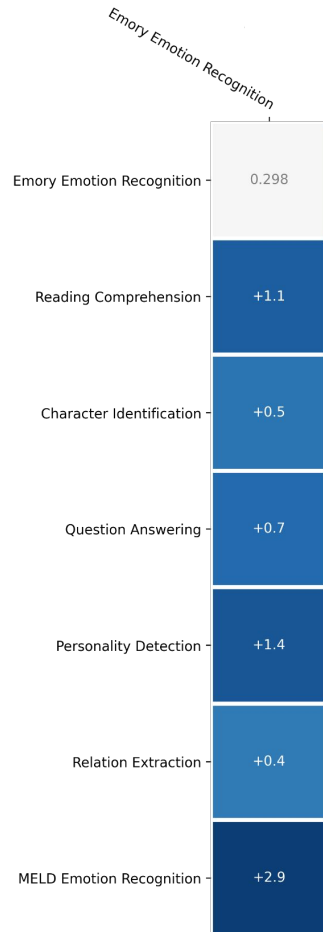
 <https://alon-albalak.github.io/>

# Supplementary Materials

# FLAD Preliminary Results

BASE MODEL AUX. DATA	T5-XL		T0-3B	
	T0MIX	P3	T0MIX	P3
Target-Only	52.82 <sub>3.34</sub>		56.44 <sub>4.70</sub>	
Explore-Only	59.18 <sub>5.52</sub>	60.64 <sub>4.92</sub>	61.17 <sub>3.30</sub>	62.77 <sub>4.83</sub>
Exploit-Only	59.79 <sub>5.63</sub>	60.49 <sub>5.01</sub>	60.87 <sub>3.35</sub>	62.87 <sub>3.69</sub>
EXP3-FLAD	61.50 <sub>4.25</sub>	64.07 <sub>4.81</sub>	62.87 <sub>3.85</sub>	<u>65.98</u> <sub>3.20</sub>
UCB1-FLAD	62.01 <sub>3.89</sub>	<u>65.52</u> <sub>3.86</sub>	62.89 <sub>3.68</sub>	<b>66.29</b> <sub>3.29</sub>

# FETA Experiments

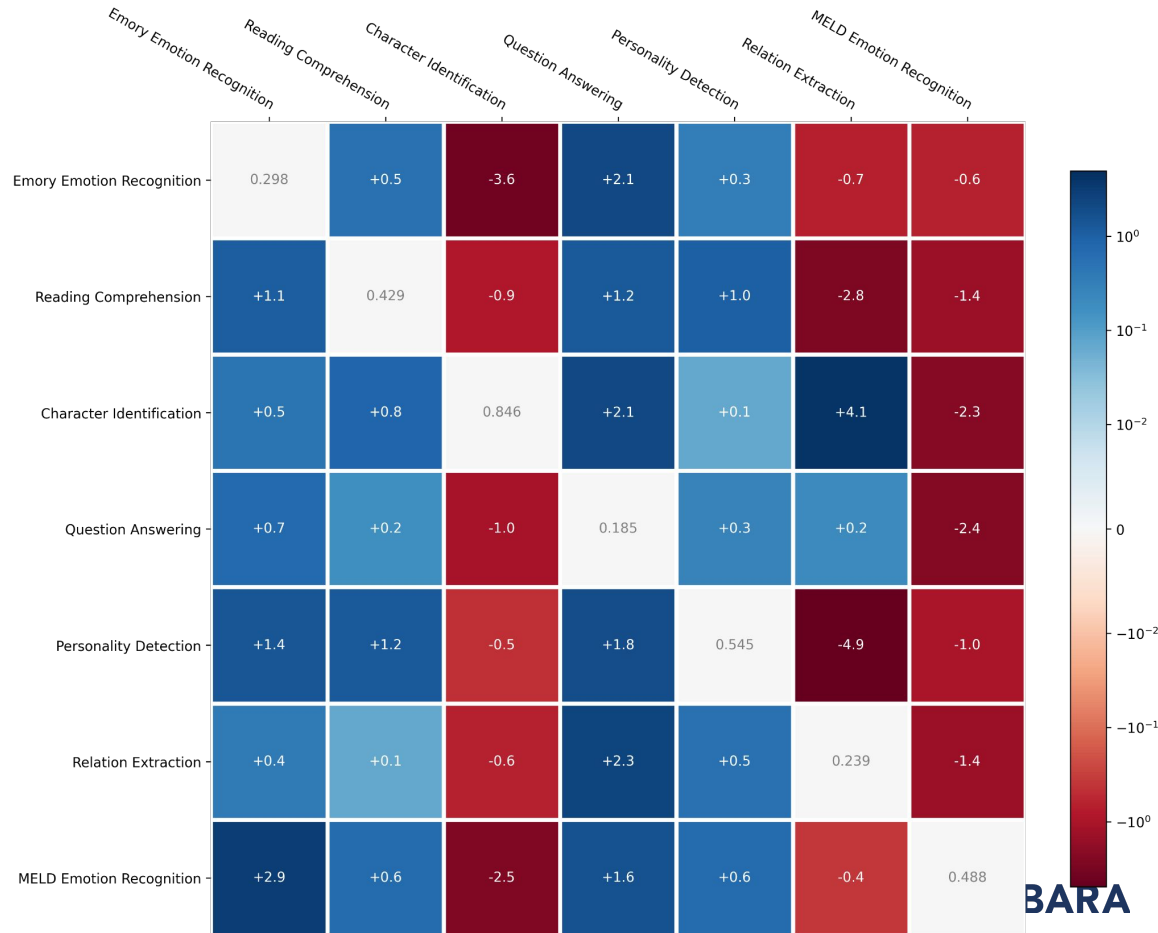


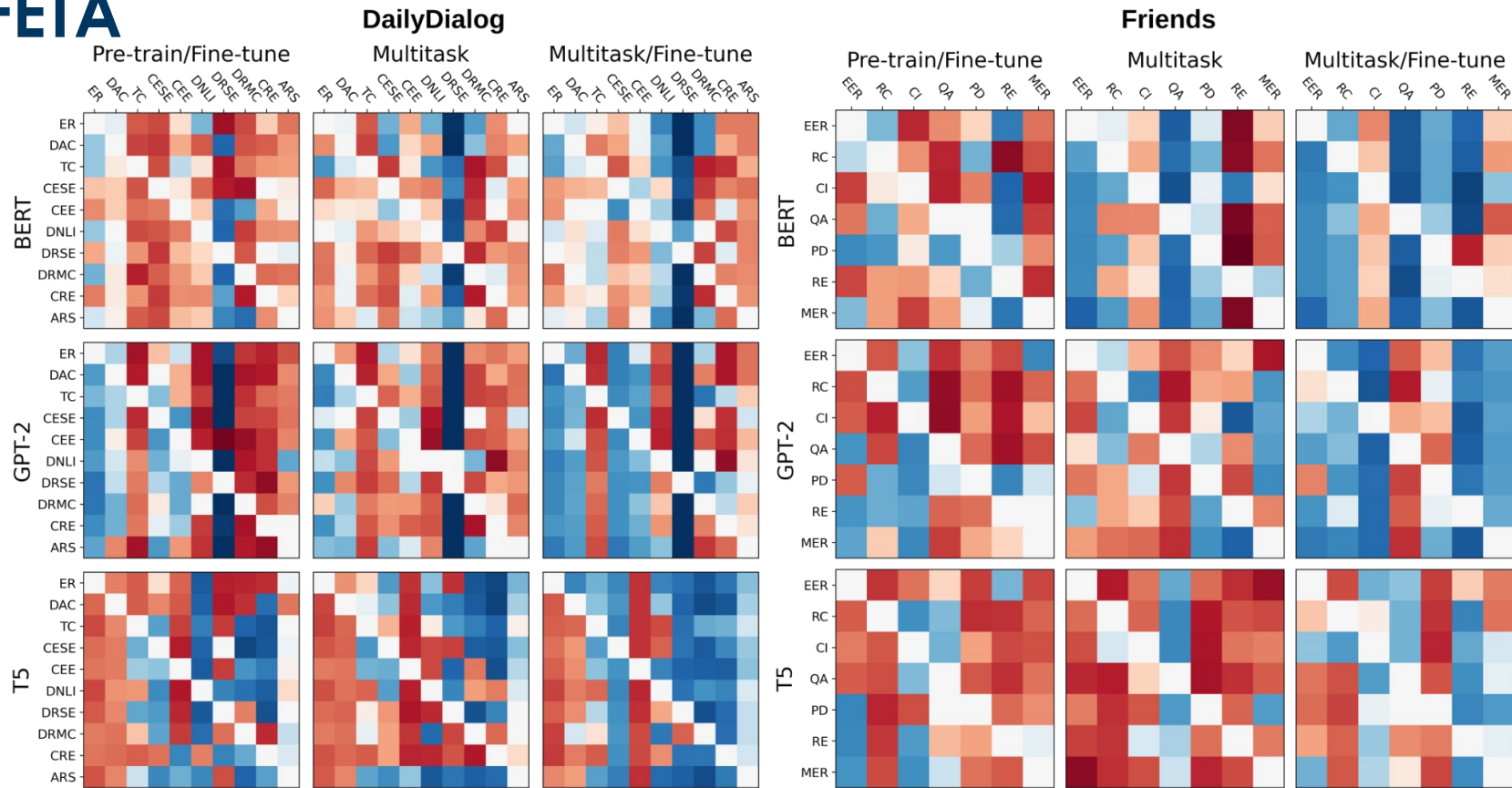
# FETA Experiments



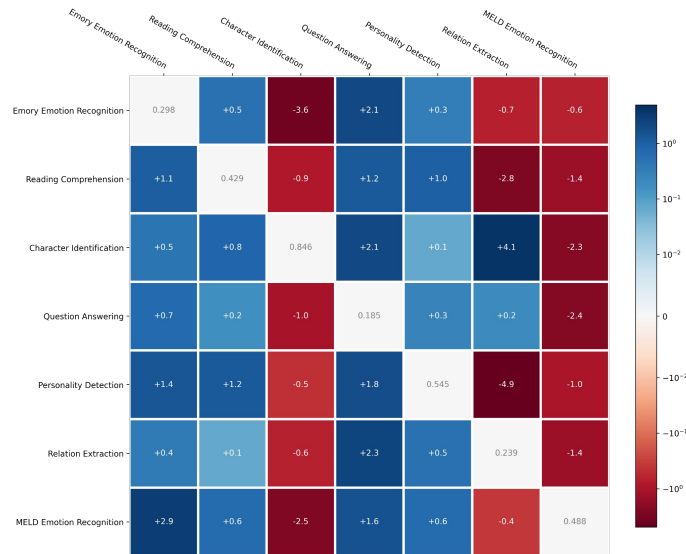


# FETA Experiments

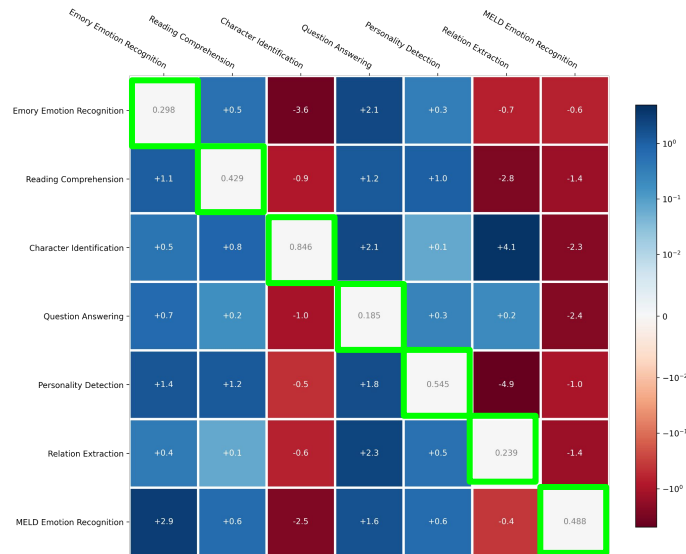




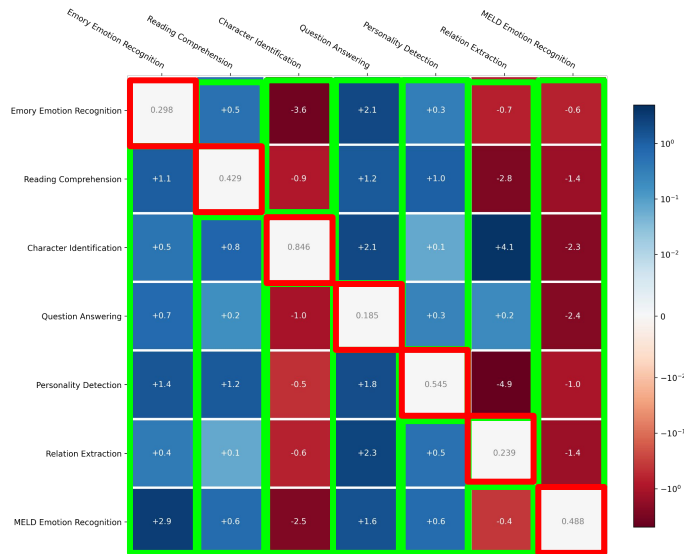
# FETA Metrics



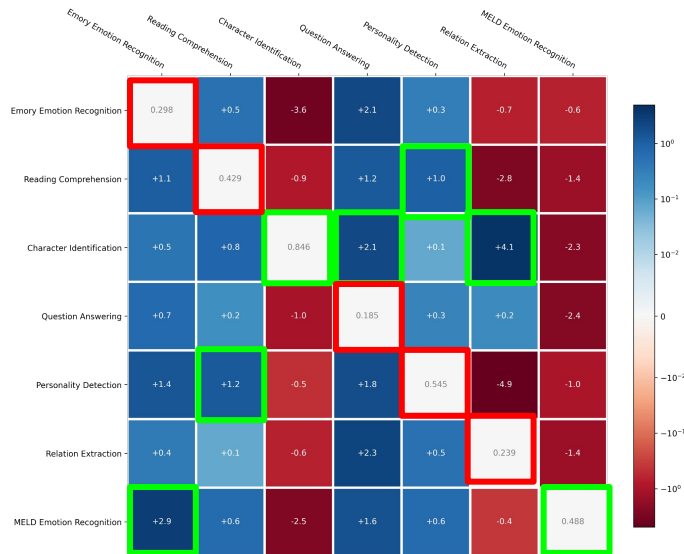
# FETA Metrics - Baseline Score



# FETA Metrics - Average Transfer Score



# FETA Metrics - Top-1 Score



# FETA Results - Aggregate Scores

Model	Transfer Algorithm	DailyDialog				Friends			
		Average		Top-1 Source		Average		Top-1 Source	
		Score ( $\sigma$ )	$\Delta$	Score	$\Delta$	Score ( $\sigma$ )	$\Delta$	Score	$\Delta$
BERT	Pre-train/Fine-tune	50.61 (0.24)	-0.93	52.22	+0.68	42.39 (0.30)	-0.89	44.36	+1.08
	Multitask	50.95 (0.24)	-0.59	52.40	+0.86	42.88 (0.29)	-0.40	45.14	+1.86
	Multitask/Fine-tune	<b>51.40</b> (0.25)	<b>-0.15</b>	<b>52.76</b>	<b>+1.22</b>	<b>44.69</b> (0.28)	<b>+1.41</b>	<b>46.00</b>	<b>+2.72</b>
GPT-2	Pre-train/Fine-tune	39.80 (0.25)	-1.28	42.19	+1.11	32.66 (0.18)	-0.64	34.34	+1.04
	Multitask	40.21 (0.24)	-0.86	41.77	+0.69	33.10 (0.16)	-0.20	34.83	+1.53
	Multitask/Fine-tune	<u>41.15</u> (0.23)	<u>+0.07</u>	<u>42.76</u>	<u>+1.68</u>	<u>34.62</u> (0.15)	<u>+1.32</u>	<u>35.86</u>	<u>+2.56</u>
T5	Pre-train/Fine-tune	49.92 (0.37)	+0.19	<b>53.04</b>	<b>+3.31</b>	41.73 (0.19)	-1.10	43.52	+0.69
	Multitask	49.49 (0.42)	-0.24	52.98	+3.25	40.42 (0.20)	-2.40	43.33	+0.51
	Multitask/Fine-tune	<u>50.29</u> (0.36)	<b>+0.56</b>	52.85	+3.12	<u>42.29</u> (0.17)	<b>-0.53</b>	<u>43.87</u>	<u>+1.05</u>

# Abstract

- Machine learning algorithms are data hungry
- Unlabeled data is easy to gather, but hard to utilize for specific tasks
- Labeled data can be time consuming and expensive to gather, and sometimes impossible due to privacy concerns or the nature of the problem
- For these reasons, few- and zero-shot settings (which require little to no labeled data) are attractive learning paradigms
- In this talk, I discuss methods of improving data efficiency in natural language processing inspired by transfer learning, reinforcement learning, and neuro-symbolic methods.
- In the end, I'll discuss current and future work for my PhD

Machine learning algorithms are data hungry, and although unlabeled data (e.g. web text) is easy to gather, it is difficult to utilize for specific tasks. Furthermore, labeled data can be time consuming and expensive to gather, and sometimes impossible due to privacy concerns or the nature of the problem.

This talk discusses methods of improving data efficiency in natural language processing inspired by transfer learning, reinforcement learning, and neuro-symbolic methods. The focus is on few- and zero-shot settings, which are attractive learning paradigms due to the challenges of gathering labeled data for specific tasks. The talk will also cover current and future work for the my PhD research in this area.



# Outline

- ~5 introductory slides
- Outline of my research area
  - What is the broad area
  - What are the specific areas I have focused on
- In depth-ish on 1-2 works
  - Multi-task learning with prompts, instruction tuning, for zero- and few-shot
  - FETA as a good testbed
  - Finish with outcome that more source datasets didn't always help, how can we improve on this?
- Future work - addressing those weaknesses + other areas of interest
  - MAB for FLAD
  -
- Timeline

# My research area

How do we improve data efficiency? Utilize additional information

We (roughly) categorize methods by the source of the information  
(2-sided figure with works on each side)

- Human knowledge (1-2 slides on how this works)
- Data related to our target domain/task

This presentation will focus on the second category

Specifically, we focus on methods for NLP

# Zero-shot transfer methods

- Background:
  - How can multiple tasks be handled by one model?
  - Encoder-only + task-specific head (classification tasks only)
  - Generative models + text-to-text format
- Multi-task learning formulation
- Instruction tuning

Benefits have been demonstrated in large and massive language models, what about small models? (slide showing benefit from previous studies on large models)

Follow ENSLP slides to discuss data/models

Transition to FETA: Next, we'll take a closer look at the effects of individual source tasks on the target task

# FETA

Transition to current/future work: So, we've seen the effects of negative transfer, not all source tasks are equal, and also that simply adding more source tasks doesn't always improve performance. So, moving forward, how can we mitigate negative transfer?