

# Exploring the Role of Task Transferability in Large-Scale Multi-Task Learning

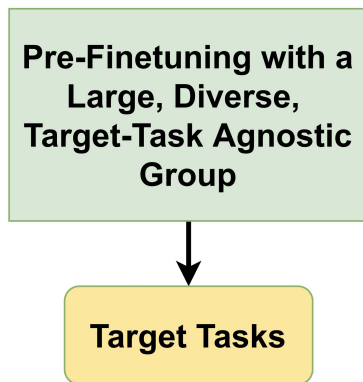
Vishakh Padmakumar<sup>1</sup>, Leonard Lausen<sup>2</sup>, Miguel Ballesteros<sup>3</sup>, Sheng Zha<sup>2</sup>, He He<sup>1,2</sup>, George Karypis<sup>2</sup>

<sup>1</sup>New York University, <sup>2</sup>AWS AI, <sup>3</sup>AWS AI Labs



# To Scale or Not To Scale, That Is The Question

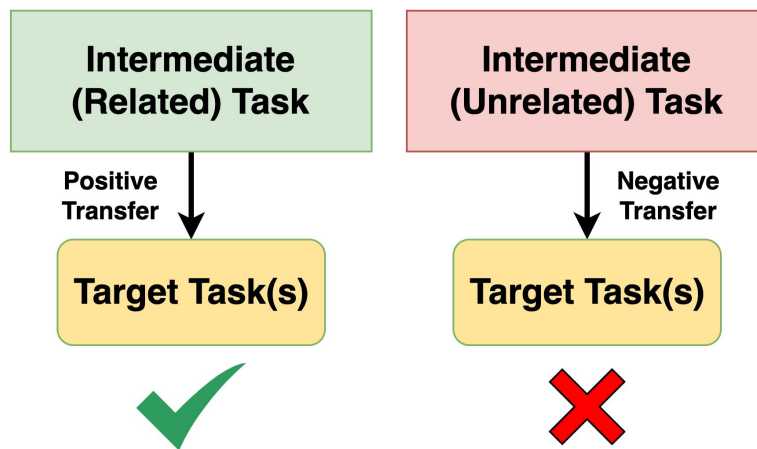
- Multi-task pre-finetuning on a **sufficiently large, diverse** set of tasks is an effective **task-agnostic** second-stage of model pre-training<sup>[1]</sup>



[1] Aghajanyan, Armen, et al. "Muppet: Massive multi-task representations with pre-finetuning." *arXiv preprint arXiv:2101.11038* (2021).

# To Scale or Not To Scale, That Is The Question

- Work on transferability has shown that the ***choice of intermediate task*** significantly impacts downstream task performance<sup>[2,3]</sup>

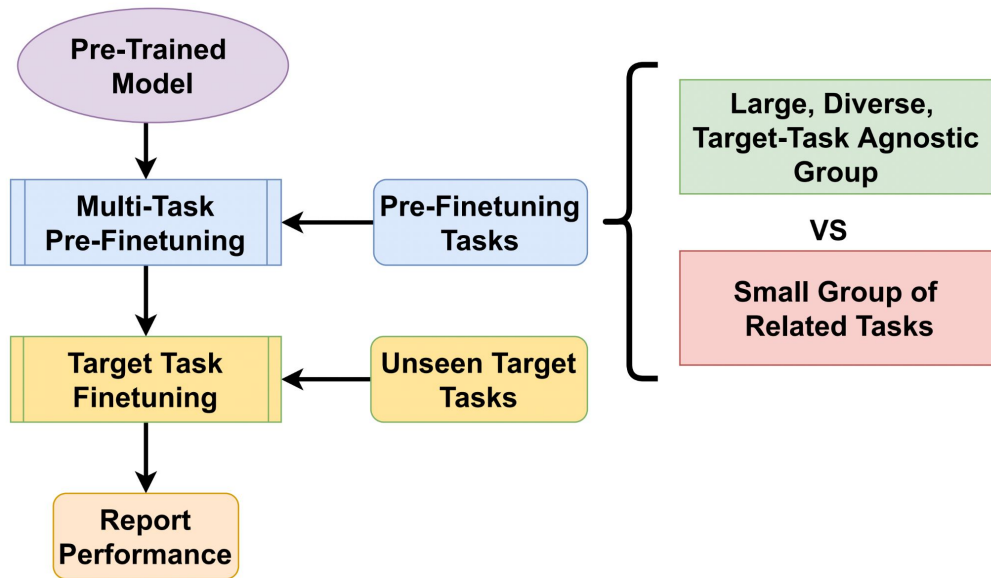


[2] Vu, Tu, et al. "Exploring and predicting transferability across NLP tasks." *arXiv preprint arXiv:2005.00770* (2020).

[3] Pruksachatkun, Yada, et al. "Intermediate-task transfer learning with pretrained models for natural language understanding: When and why does it work?." *arXiv preprint arXiv:2005.00628* (2020).

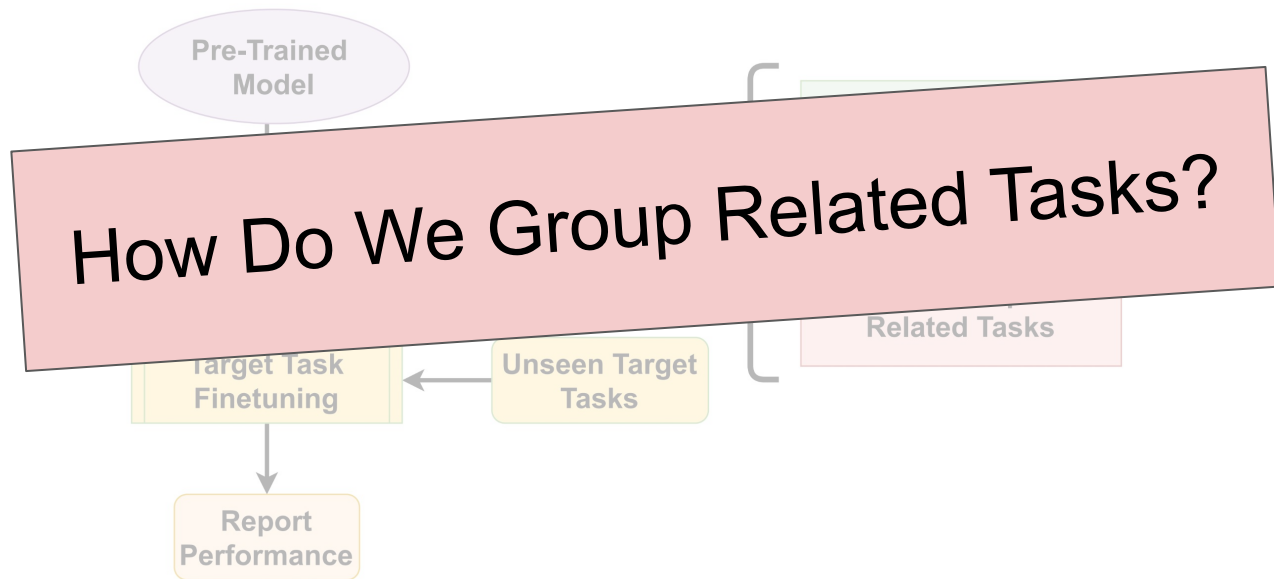
# Research Question

- We aim to study how the **choice of pre-finetuning tasks** and the **size of the multi-task step** affects target task performance.

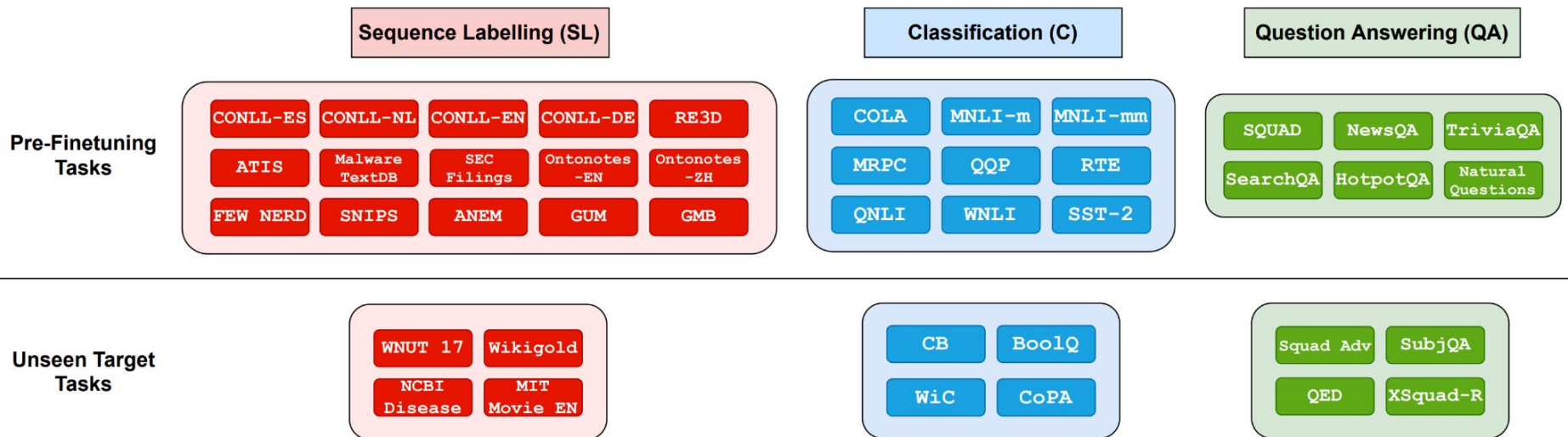


# Research Question

- We aim to study how the *choice of pre-finetuning tasks* and the *size of the multi-task step* affects target task performance.



# How Do We Group Related Tasks?



# Experiments

Run pre-finetuning on each combination of task groups

	<b>Baseline</b>	<b>Only-SL</b>	<b>Only-C</b>	<b>Only-QA</b>	<b>SL+C</b>	<b>SL+QA</b>	<b>QA+C</b>	<b>SL+C+QA</b>
<b>Unseen SL</b>	80.134	<b>80.844</b>	72.370	78.095	79.455	80.105	77.378	80.750
<b>Unseen C</b>	68.109	67.406	71.404	63.21	70.422	70.068	70.067	<b>73.021</b>
<b>Unseen QA</b>	56.692	45.174	61.120	75.252	57.568	75.460	75.035	<b>75.678</b>
<b>Average</b>	68.312	64.475	68.298	72.385	69.147	75.564	74.160	<b>76.483</b>

# Experiments

Report results on all unseen tasks, averaged over the task groups

	Baseline	Only-SL	Only-C	Only-QA	SL+C	SL+QA	QA+C	SL+C+QA
<b>Unseen SL</b>	80.134	<b>80.844</b>	72.370	78.693	79.453	80.165	77.378	80.750
<b>Unseen C</b>	68.109	67.406	71.404	63.21	70.422	70.068	70.067	<b>73.021</b>
<b>Unseen QA</b>	56.692	45.174	61.120	75.252	57.568	75.460	75.035	<b>75.678</b>
<b>Average</b>	68.312	64.475	68.298	72.385	69.147	75.564	74.160	<b>76.483</b>



# Effect of Multi-Task Scaling

On average, a large-scale task-agnostic multi-task step improves downstream performance

	Baseline	Only-SL	Only-C	Only-QA	SL+C	SL+QA	QA+C	SL+C+QA
Unseen SL	80.134	<b>80.844</b>	72.370	78.693	79.453	80.165	77.378	80.750
Unseen C	68.109	67.406	71.404	63.21	70.422	70.068	70.068	<b>73.021</b>
Unseen QA	56.692	45.174	61.120	75.252	57.568	75.460	75.038	<b>75.678</b>
Average	68.312	64.475	68.298	72.385	69.147	75.564	74.160	<b>76.483</b>

# Effect of Transferability

Target task performance on unseen tasks is improved when we pre-finetuning on related tasks from the same group

	Baseline	Only-SL	Only-C	Only-QA	SL+C	SL+QA	QA+C	SL+C+QA
<b>Unseen SL</b>	80.134	<u>80.844</u>	<u>72.370</u>	<u>78.693</u>	79.453	80.165	77.378	80.750
<b>Unseen C</b>	68.109	<u>67.406</u>	<u>71.404</u>	<u>63.21</u>	70.422	70.068	70.067	<b>73.021</b>
<b>Unseen QA</b>	56.692	<u>45.174</u>	<u>61.120</u>	<u>75.252</u>	57.568	75.460	75.035	<b>75.678</b>
<b>Average</b>	68.312	64.475	68.298	72.385	69.147	75.564	74.160	<b>76.483</b>

# Interplay of Transferability and Multi-Task Scaling

Target task performance on pre-finetuning with a small group of related tasks is on-par with the large-scale multi-task setup

	Baseline	Only-SL	Only-C	Only-QA	SL+C	SL+QA	QA+C	SL+C+QA
<b>Unseen SL</b>	80.134	<u>80.844</u>	72.370	78.693	79.453	80.165	77.378	<u>80.750</u>
<b>Unseen C</b>	68.109	67.406	<u>71.404</u>	63.21	70.422	70.068	70.067	<u>73.021</u>
<b>Unseen QA</b>	56.692	45.174	61.120	<u>75.252</u>	57.568	75.460	75.035	<u>75.678</u>
<b>Average</b>	68.312	64.475	68.298	72.385	69.147	75.564	74.160	<b>76.483</b>

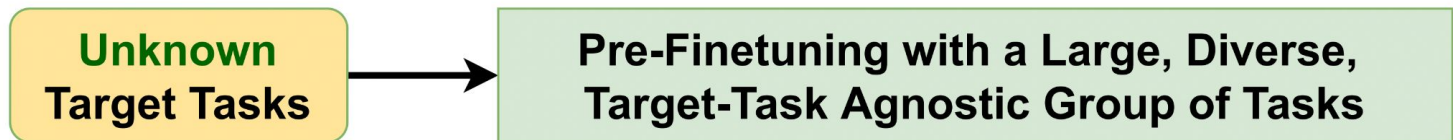
# Interplay Between Task Groups

It's hard to predict the interplay between tasks, so selecting an optimum subset of tasks is challenging

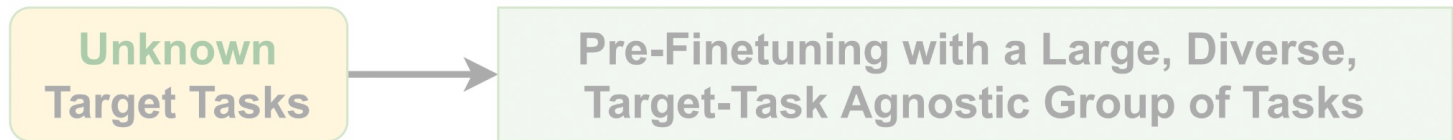
	Baseline	Only-SL	Only-C	Only-QA	SL+C	SL+QA	QA+C	SL+C+QA
<b>Unseen SL</b>	80.134	<b>80.844</b>	72.370	78.693	<u>79.453</u>	<u>80.165</u>	<u>77.378</u>	80.750
<b>Unseen C</b>	68.109	67.406	71.404	63.21	<u>70.422</u>	<u>70.068</u>	<u>70.067</u>	<b>73.021</b>
<b>Unseen QA</b>	56.692	45.174	61.120	75.252	<u>57.568</u>	<u>75.460</u>	<u>75.035</u>	<b>75.678</b>
<b>Average</b>	68.312	64.475	68.298	72.385	69.147	75.564	74.160	<b>76.483</b>

# Takeaways

- When the target tasks are **unknown**, **multi-task scaling** provides an effective **general purpose model**



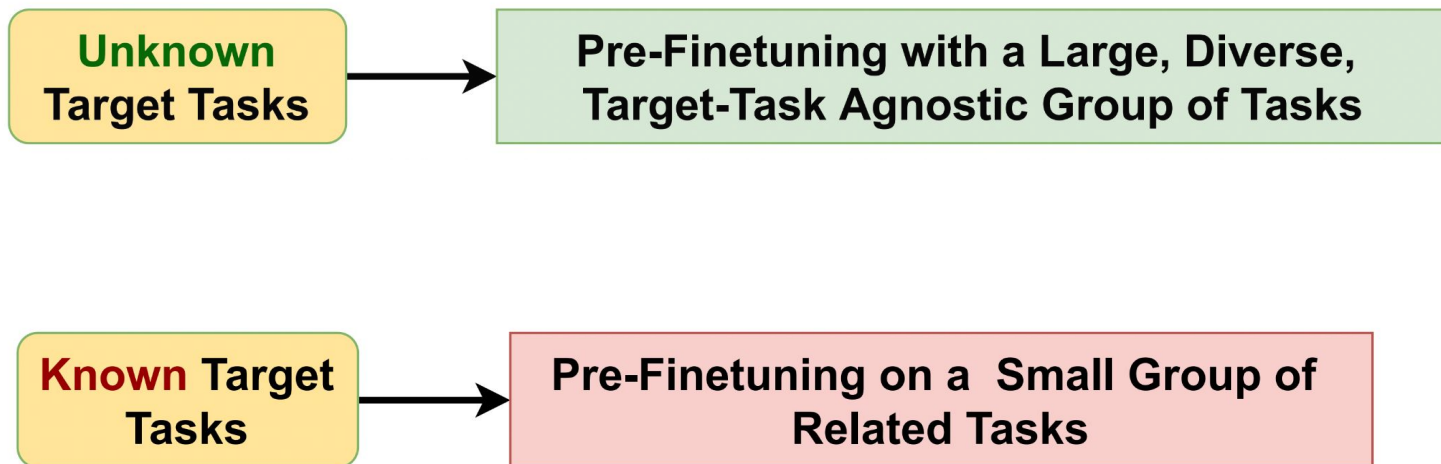
# Takeaways



- If the goal is to improve performance on a **specific target task(s)** then a **smaller set of related tasks** is an effective, **computationally cheaper alternative**



# Takeaways



For more details, stop by  
Poster Session 2 on 7/11 at 2:30pm :)