

# Modular Deep Learning

Edoardo M. Ponti  
[eponti@ed.ac.uk](mailto:eponti@ed.ac.uk)

Sapienza, 9 Jun 2023



THE UNIVERSITY  
*of* EDINBURGH



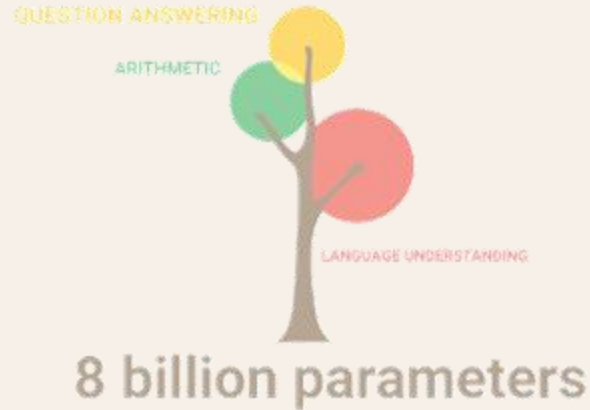
UNIVERSITY OF  
CAMBRIDGE

# Background and Motivation

Pfeiffer, Jonas, Sebastian Ruder, Ivan Vulić, and Edoardo M. Ponti.  
*Modular deep learning. arXiv 2023*



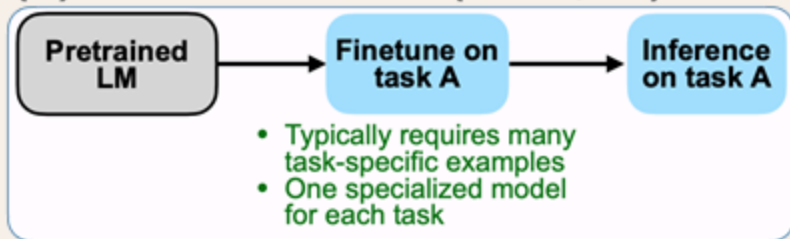
# Emergence of NLP Abilities via Scaling



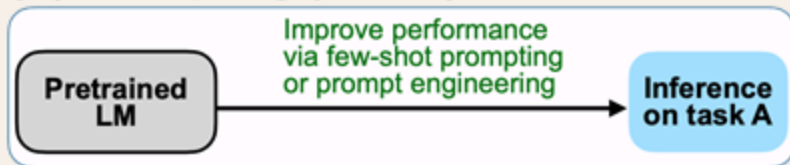
Credits: Google AI Blog

# Zero/few-shot generalisation to new tasks

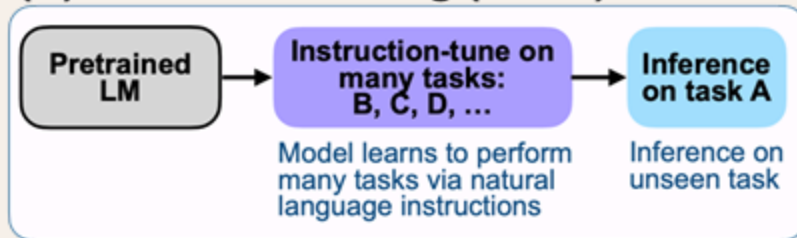
## (A) Pretrain–finetune (BERT, T5)



## (B) Prompting (GPT-3)



## (C) Instruction tuning (FLAN)

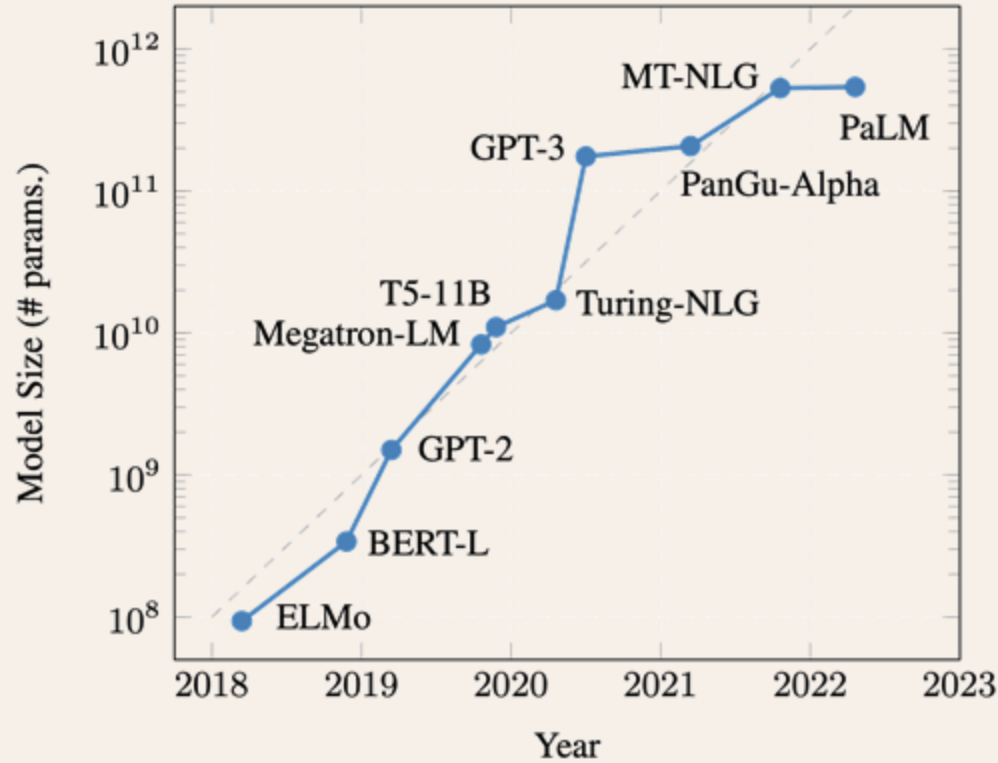


[Wei et al. 2022]

# Challenges of Existing Methods

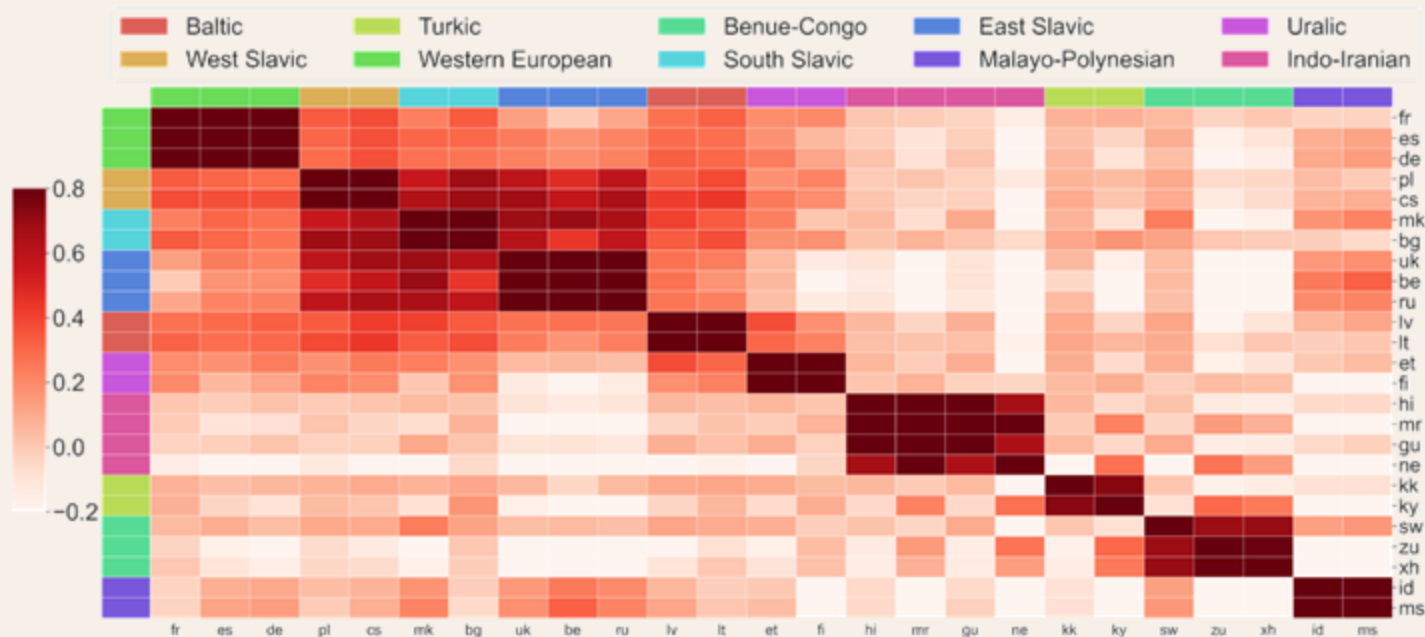
- **Inefficiency due to model size:**  
time and space complexity
- **Negative transfer:**  
interference (multi-task learning)  
catastrophic forgetting (continuous learning)
- **Systematic generalisation:**  
sub-problem recombination  
local distribution shifts

# Inefficiency due to Model Size



Evolution of the size of large pre-trained models [\[Treviso et al., 2022\]](#)

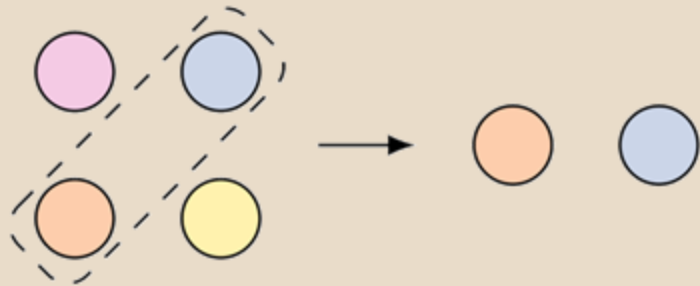
# Negative Transfer



Cosine similarity between gradients for NMT from multiple languages to English [Wang et al. 2021, Gradient Vaccine].

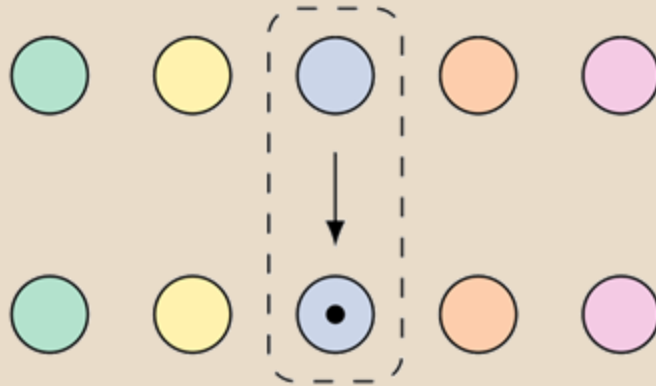
# Systematic Generalisation [Hupkes et al. 2020]

Disentangling and recombining autonomous facets of knowledge



**Skill recombination**

Adapting to local distribution shifts based on few examples



**Local adaptation**



# Modular Deep Learning [Pfeiffer et al. 2023]

Correspondence between **modules** of a network and the **specialised functions** they perform.

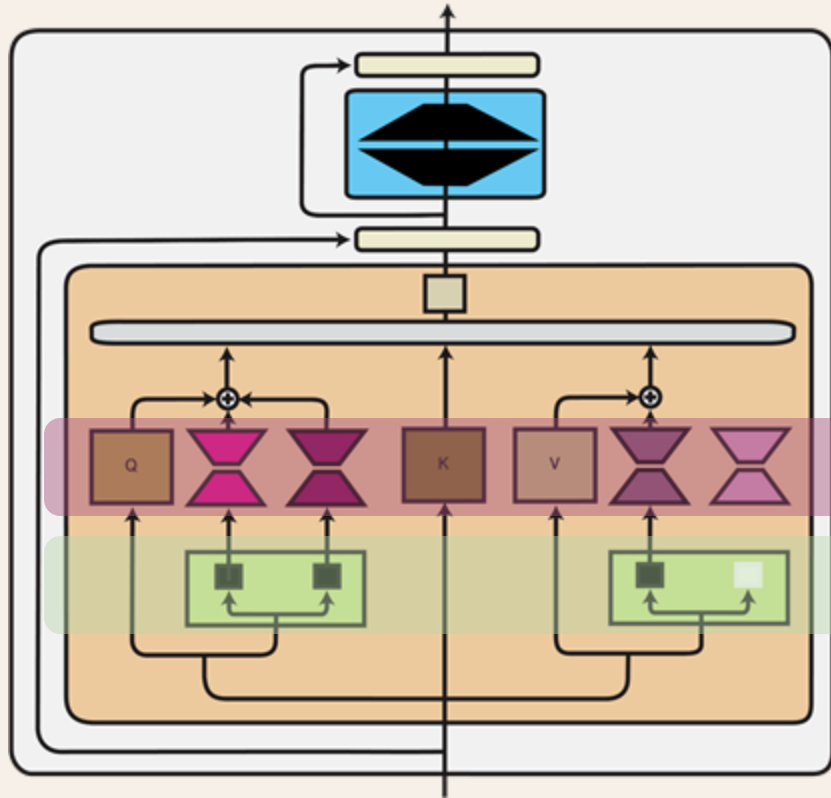
Separation between **routing** (controlling information flow) and **computation** (modules).

In **multitask learning**, modules should specialise for **skills**, subproblems common to multiple tasks which can be recombined or locally updated.

# Advantages

- **Over in-context learning:**
  - **Robustness** (no instability due to ordering [Zhao et al., 2021], wording [Webson & Pavlick, 2022], etc.)
  - **Higher performance** [Brown et al. 2020]
- **Over (instruction) fine-tuning:**
  - Only **positive transfer** across tasks (no catastrophic **forgetting and interference**) [Caccia et al. 2023]
  - Compositionality, reusability, and local updates of modules: **systematic generalisation** [Ponti et al. 2020]
  - **Parameter efficiency** (no large full-model copies) [Liu et al., 2022]
  - **Scaling** (e.g., through MoE) [Shazeer et al. 2017]

# A Blueprint of Modular Deep Learning

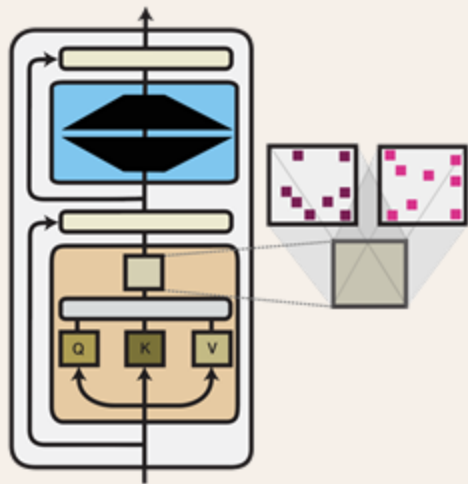


1) Modules

2) Routing

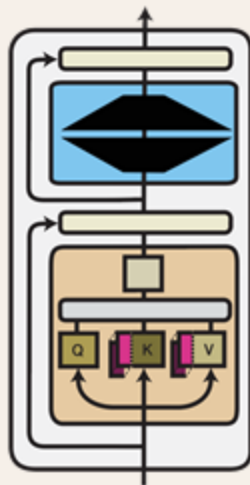
3) Applications

# Modules: Parameter-Efficient Fine-Tuning



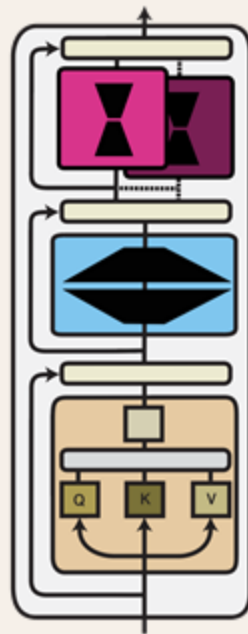
Parameter Composition

$$f'_i(\mathbf{x}) = f_{\theta_i \oplus \phi}(\mathbf{x})$$



Input Composition

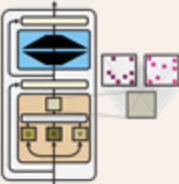
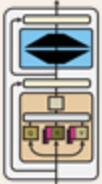
$$f'_i(\mathbf{x}) = f_{\theta_i}([\mathbf{x}, \phi])$$



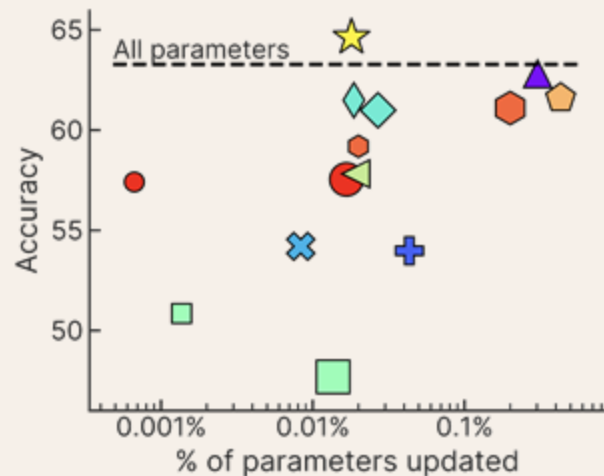
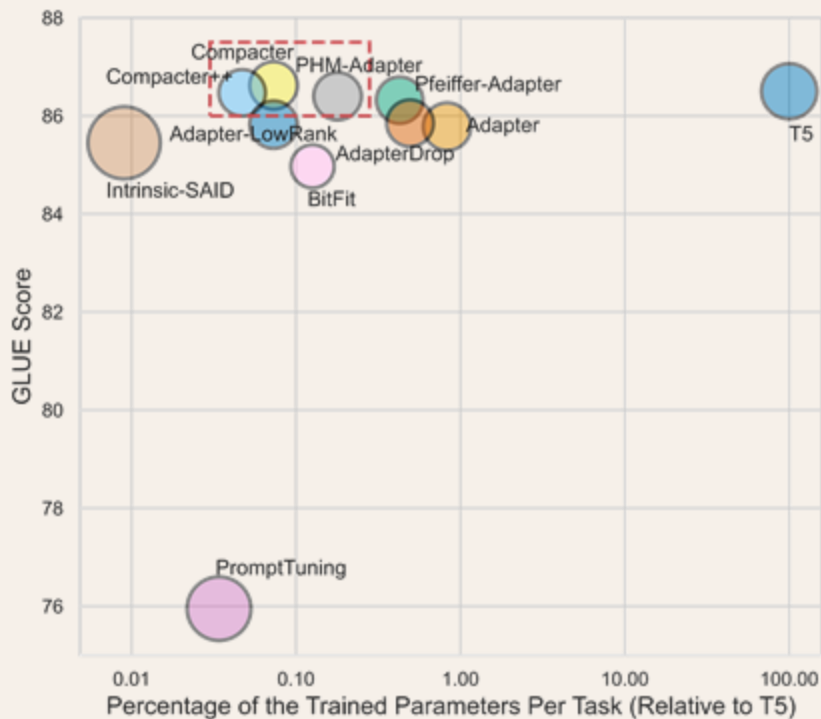
Function Composition

$$f'_i(\mathbf{x}) = f_{\theta_i} \odot f_{\phi}(\mathbf{x})$$

# Modules: Parameter-Efficient Fine-Tuning

	Parameter efficiency	Training efficiency	Inference efficiency	Performance	Compositionality
Parameter composition 	+	-	++	+	+
Input composition 	++	--	--	-	+
Function Composition 	-	-	-	++	+

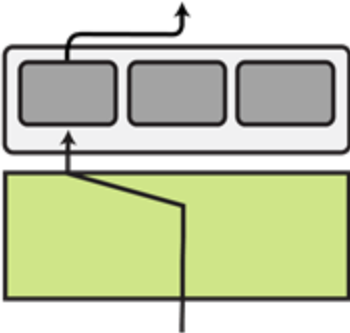
# Performance and Efficiency Comparison



Performance, param efficiency, and memory footprint of different methods on T5-Base (222Mparams; left) [\[Mahabadi et al., 2021\]](#) and T5-3B (3B params; right) [\[Liu et al., 2022\]](#)

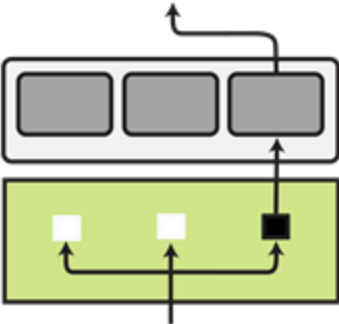
# Routing

Fixed Routing

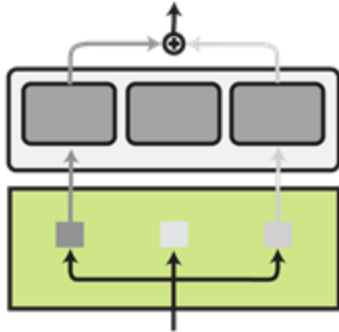


Learned Routing

Hard Learned Routing



Soft Learned Routing



# Applications

- **Zero-shot Cross-lingual Transfer**
- **Faithful and Abstractive Dialogue Generation**
- **Few-shot Adaptation to New RL / NLP Tasks**
- ...and many more (including programme induction and causal inference / discovery)!



# Survey on Modular Deep Learning

<https://www.modulardeeplearning.com>

## TALKS

*Talks related to our survey paper.*

### Modular and Parameter-Efficient Fine-Tuning for NLP Models

Sebastian Ruder, Jonas Pfeiffer, Ivan Vulčić  
EMNLP 2022, December 8, 2022



### Modular and Parameter- Efficient Fine-Tuning for NLP Models

*EMNLP 2022 Tutorial*

← Cohere For AI




Jonas Pfeiffer  
Research Scientist  
Google Research

Modular and  
Composable  
Transfer Learning

Thursday, October 6  
11AM ET

### Modular and Composable Transfer Learning

*Jonas Pfeiffer @ Cohere for AI*



### Combining modular skills in multitask learning

*Edoardo M. Ponti*

*Efficient Large-Scale AI Workshop  
Microsoft Research, 2022/10/24*

### Combining modular skills in multitask learning

*Edoardo M. Ponti @ Microsoft Research Summit  
2022*



# Composition of Sparse Adapters

Ansell, Alan, Edoardo M. Ponti, Anna Korhonen, and Ivan Vulić.

*Composable Sparse Fine-Tuning for Cross-Lingual Transfer. ACL 2022*



# Standard Zero-Shot Cross-Lingual Transfer

Step 1:

**Pre-train** a multilingual model.



MODEL

Step 2:

**Fine-tune** the model on a **task** in a high-resource **source language**.

Step 3:

Transfer and **evaluate** the model on a low-resource **target language**.

**Why?**

Training **data** is **expensive** and not available for many languages, especially ones that are considered “low-resource”.

# Modular Zero-Shot Cross-Lingual Transfer

## Step 1: Train Language Adapters

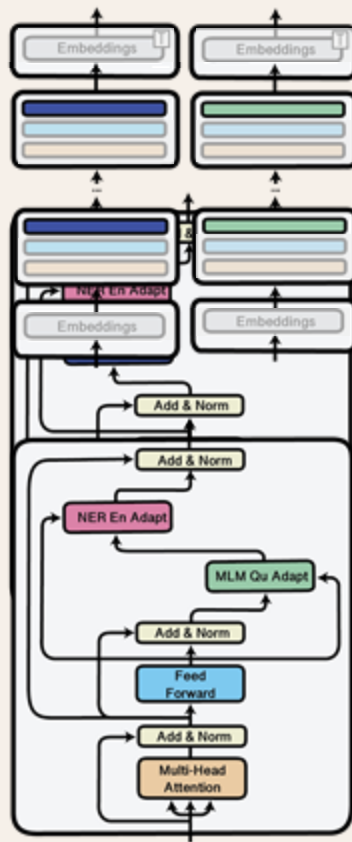
We train **language adapters** for the **source language** and the **target language** with masked language modelling on Wikipedia.

## Step 2: Train a Task Adapter

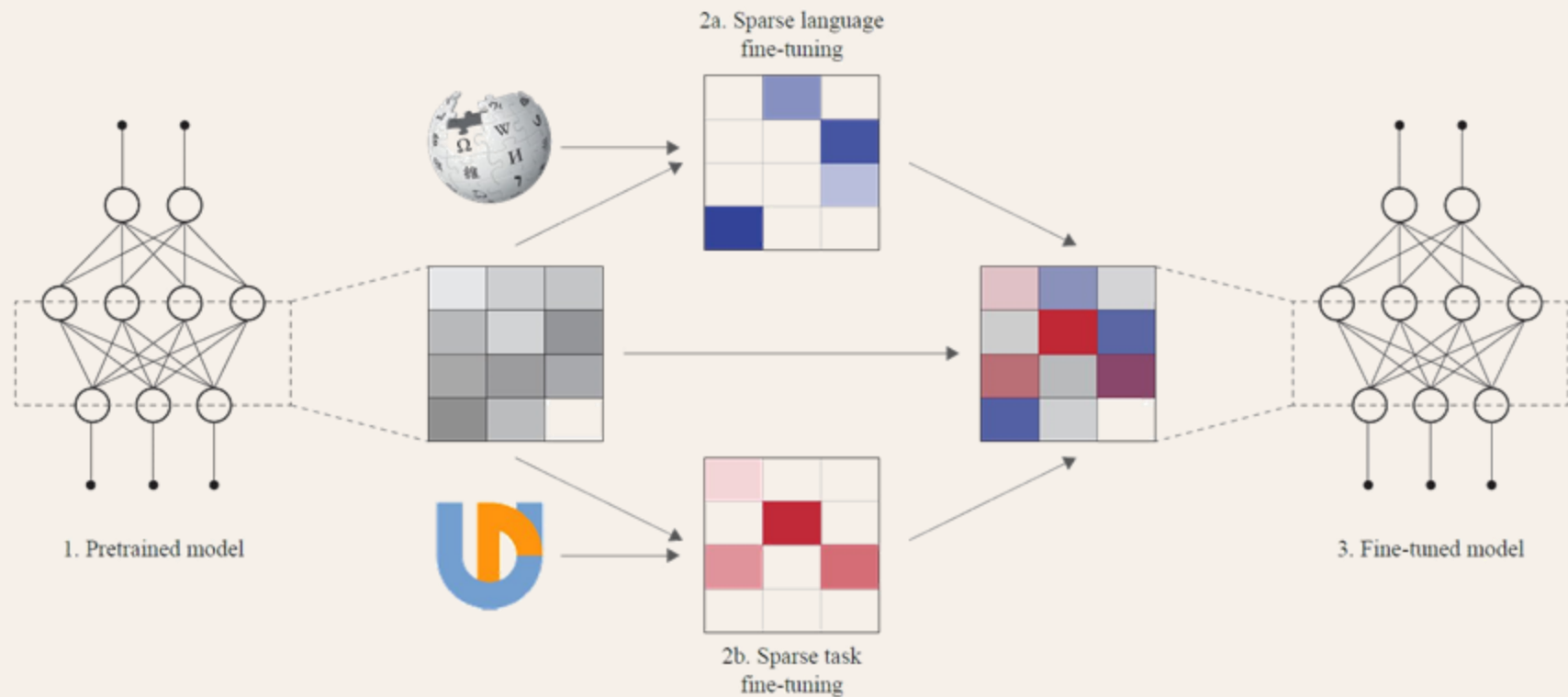
We **train task adapters** in the source language **stacked** on top of the source **language adapter**.

## Step 3: Zero-Shot transfer to unseen language

We **replace** the **source language adapter** with the **target language adapter**, while **keeping** the “language agnostic” **task adapter**.



# Sparse Fine-Tuning



# Lottery Ticket-inspired Algorithm

**Algorithm 1** Cross-Lingual Transfer with Lottery-Ticket Sparse Fine-Tuning

**function** LTSFT( $\mathcal{D}$ ,  $\mathcal{L}$ ,  $\theta^{(0)}$ ,  $\eta$ ,  $K$ )

$\theta^{(1)} \leftarrow \theta^{(0)}$

**while** not converged **do**

$\theta^{(1)} \leftarrow \theta^{(1)} - \eta \nabla \mathcal{L}(\theta^{(1)}, \mathcal{D})$

$\mu_i \leftarrow \begin{cases} 1 & \text{if } \theta_i^{(1)} \in \operatorname{argmax}_{\theta_1, \dots, \theta_K} |\theta^{(1)} - \theta^{(0)}| \\ 0 & \text{otherwise} \end{cases}$

$\theta^{(2)} \leftarrow \theta^{(0)}$

**while** not converged **do**

$\theta^{(2)} \leftarrow \theta^{(2)} - \mu \odot \eta \nabla \mathcal{L}(\theta^{(2)}, \mathcal{D})$

$\phi \leftarrow \theta^{(2)} - \theta^{(0)}$

**return**  $\phi$

**end function**

**function** CROSSLINGUALTRANSFER( $\mathcal{D}_{\text{src}}$ ,  $\mathcal{D}_{\text{tar}}$ ,  $\mathcal{D}_{\text{task}}$ ,  $\mathcal{L}_{\text{task}}$ ,  $\theta^{(0)}$ ,  $\eta$ ,  $K$ )

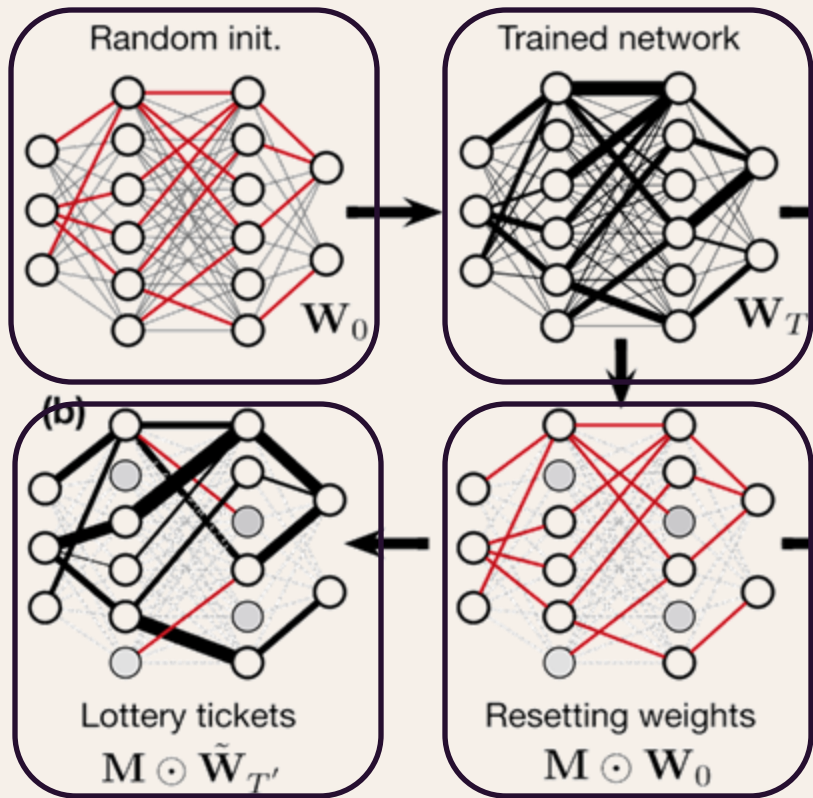
$\phi_{\text{src}} \leftarrow \text{LTSFT}(\mathcal{D}_{\text{src}}, \mathcal{L}_{\text{MLM}}, \theta^{(0)}, \eta, K)$

$\phi_{\text{task}} \leftarrow \text{LTSFT}(\mathcal{D}_{\text{task}}, \mathcal{L}_{\text{task}}, \theta^{(0)} + \phi_{\text{src}}, \eta, K)$

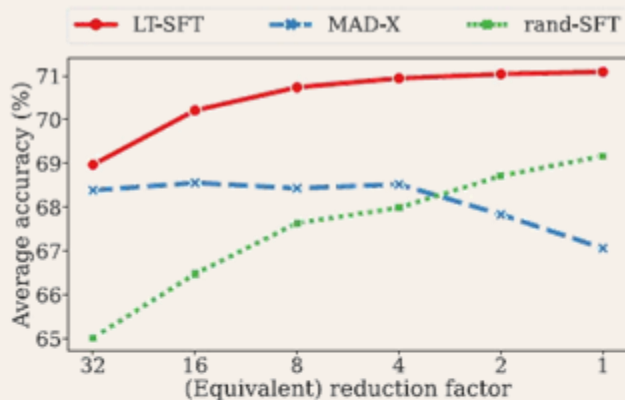
$\phi_{\text{tar}} \leftarrow \text{LTSFT}(\mathcal{D}_{\text{tar}}, \mathcal{L}_{\text{MLM}}, \theta^{(0)}, \eta, K)$

**return**  $\theta^{(0)} + \phi_{\text{task}} + \phi_{\text{tar}}$

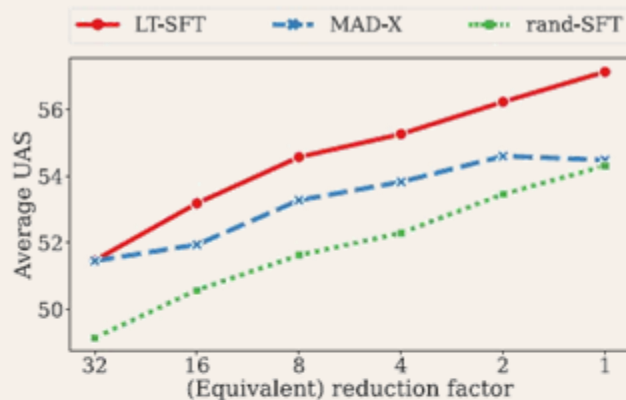
**end function**



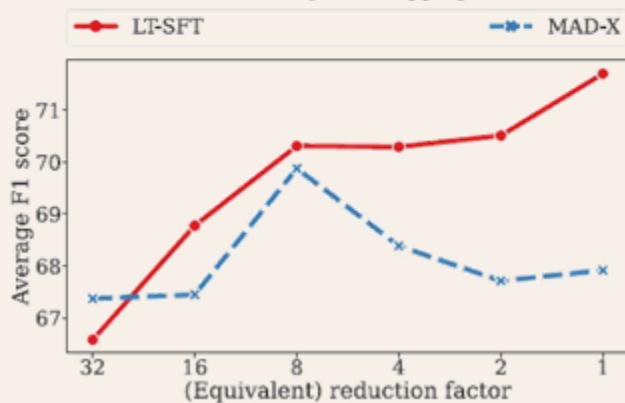
# Results for Zero-shot Transfer



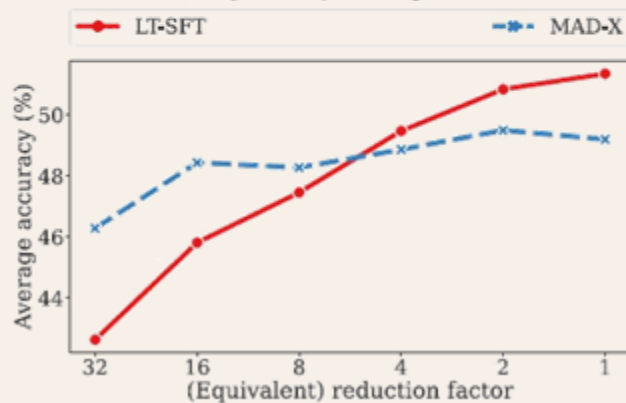
(a) Part-of-Speech Tagging



(b) Dependency Parsing (DP)



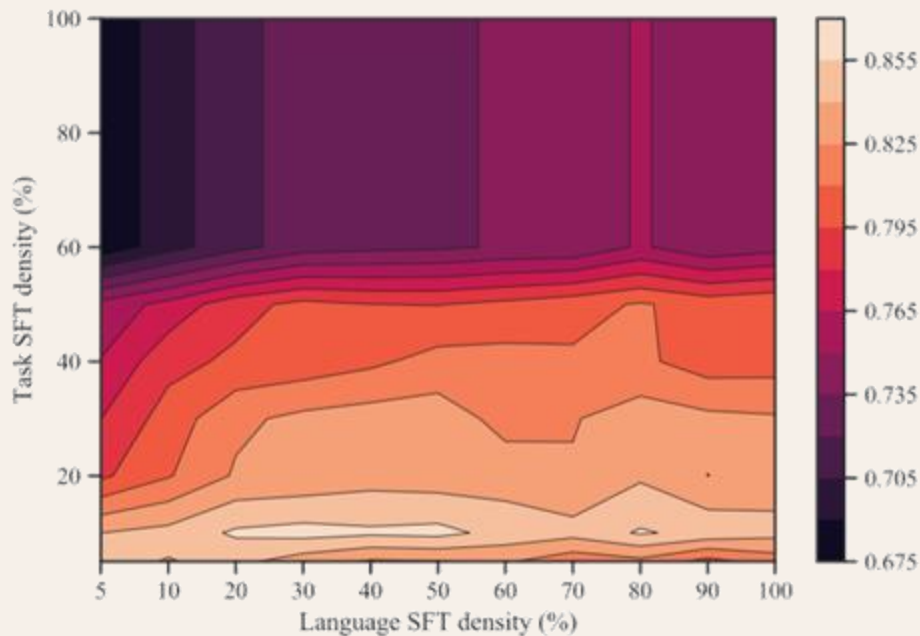
(c) Named Entity Recognition (NER)



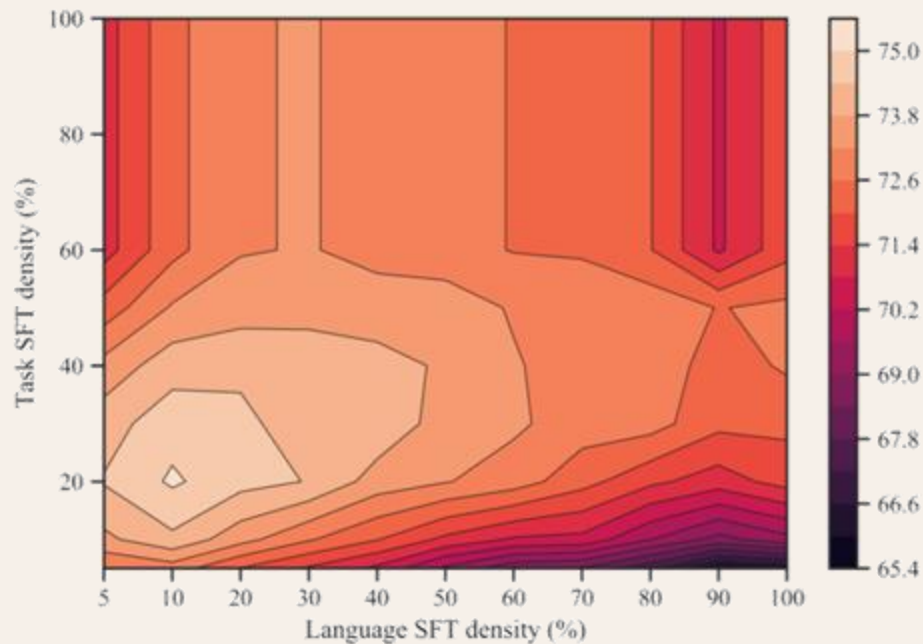
(d) Natural Language Inference (NLI)

# On the Importance of Sparsity

F1 score in Hausa NER



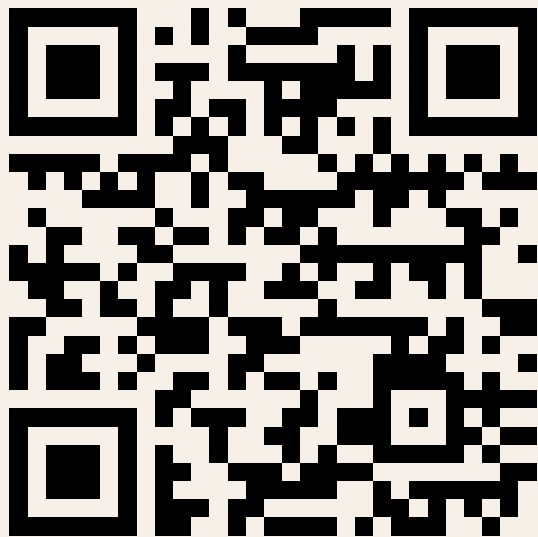
UAS score in Maltese DP





# Code

[github.com/  
cambridge-tl/  
composable-sft](https://github.com/cambridge-tl/composable-sft)



## Language SFTs

Identifiers for language SFTs are of the form `cambridge-tl/{base_model}-lang-sft-{lang_code}-small`, e.g. `cambridge-tl/mbert-lang-sft-en-small`. "Small" SFTs have ~7.6M parameters - we may release larger models in the future. Language SFTs are currently available for the following languages/models:

Language	Code	bert-base-multilingual-cased (mbert)	xlm-roberta-base (xlmr)
Acehnese	ace	×	✓
Amharic	amh	×	✓
Arabic	ar	✓	✓
Ashaninka	cni	×	✓
Balinese	ban	×	✓
Bambara	bm	✓	×
Banjarese	bjn	×	✓
Basque	eu	✓	×
Bengali	bn	✓	×
Bribri	bzd	×	✓
Bulgarian	bg	×	✓
Buryat	bxr	✓	×
Cantonese	yue	✓	×
Chinese	zh	✓	✓
Czech	cs	✓	×
English	en	✓	✓
Erzya	myv	✓	×
Estonian	et	✓	×
Faroese	fo	✓	×
French	fr	✓	✓
German	de	✓	✓

# Faithful and Abstractive Dialogue Generation

Daheim, Nico, Nouha Dziri, Mrinmaya Sachan, Iryna Gurevych, and Edoardo M. Ponti. *Elastic Weight Removal for Faithful and Abstractive Dialogue Generation*. arXiv 2023



# Knowledge-grounded dialogue generation

$$p_{\theta}(u_{T+1} \mid u_1^T, \hat{\mathcal{K}}) = \prod_{n=1}^{N_{T+1}} p_{\theta}([u_{T+1}]_n \mid [u_{T+1}]_1^{n-1}, u_1^T, \hat{\mathcal{K}})$$

Faithfulness (opposite: **hallucination**)  
is the adherence of the generated  
response  $u$  to the knowledge  $K$

Abstractive (opposite: **extractive**)  
responses  $u$  do not copy-paste spans  
from knowledge  $K$  but rephrase them.

$\mathcal{K}$ : The Flash first appeared in “Showcase” #4 (October 1956) [...]

$u_T$ : What comic series is he from?

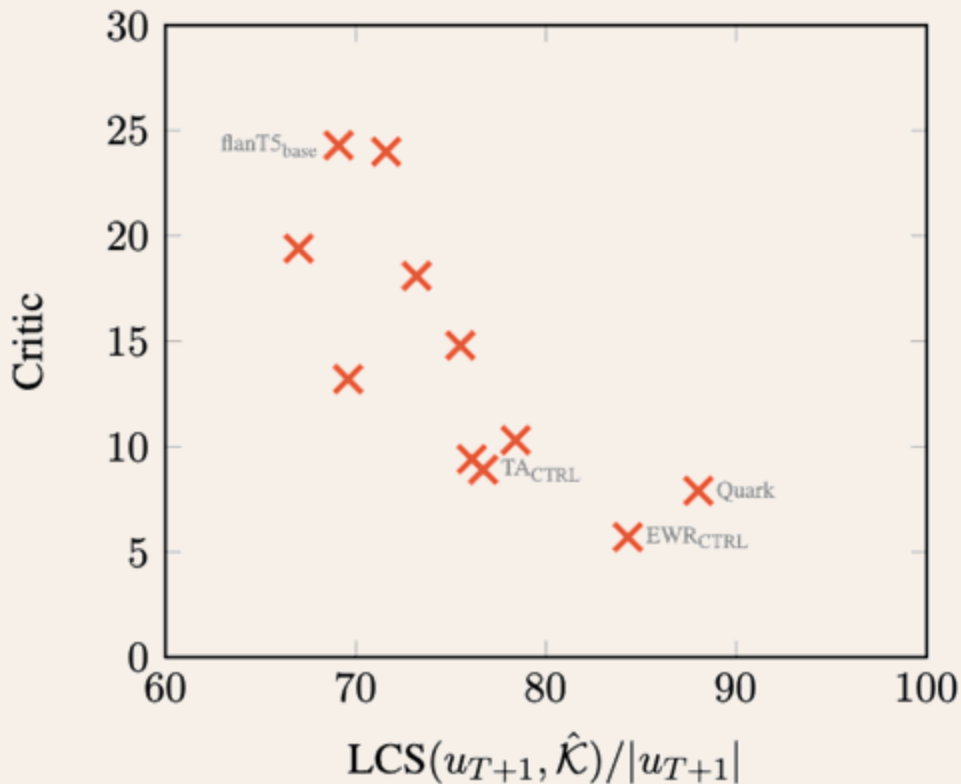
$u_{T+1}$	F	A
He first appeared in “Showcase” #4 (November 1956).	✗	✗
He first appeared in “Showcase” #4 (October 1956).	✓	✗
His first appearance was in Showcase #4 in October 1956.	✓	✓

# Faithfulness—Abtractiveness Trade-off

Faithfulness can be measured by a Critic (a binary classifier)

Abtractiveness can be measured by normalised LCS (longest common span)

Faithful models generally incur **extractive** generation. Can we have the best of both worlds?



# Task Arithmetic

[Ilharco et al. 22]

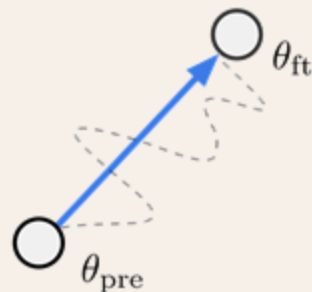
Step 1: Create **task vectors** as the difference between a model fine-tuned on examples of (positive / negative) behaviour and a pre-trained model.

$$\tau_i \triangleq \theta_i - \theta_0$$

Step 2: Add / Subtract task vectors to the pre-trained model.

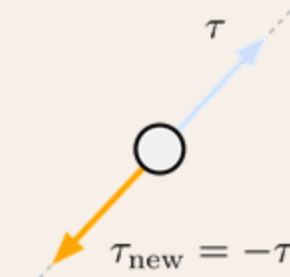
$$\theta' = \theta_0 + \sum_i \lambda_i \tau_i$$

a) Task vectors



$$\tau = \theta_{ft} - \theta_{pre}$$

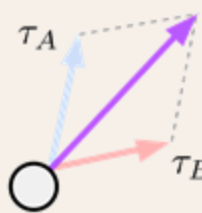
b) Forgetting via negation



Example: making a language model produce less toxic content

c) Learning via addition

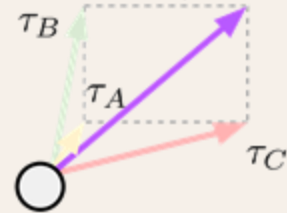
$$\tau_{new} = \tau_A + \tau_B$$



Example: building a multi-task model

d) Task analogies

$$\tau_{new} = \tau_C + (\tau_B - \tau_A)$$



Example: improving domain generalization

# Elastic Weight Addition and Subtraction

Limitations of task arithmetic:

- 1) task vectors may **interfere** with each other;
- 2) individual parameters have **higher importance** than others in controlling a certain behaviour;

Solution: weight task vectors by **Fisher Information**  $f$  to prevent interference and represent parameter importance.

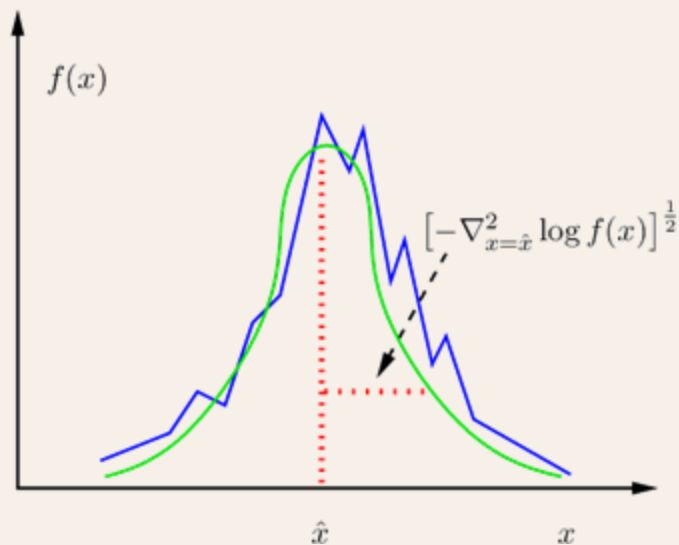
$$\frac{\lambda_0 \cdot f_{\theta_0} \cdot \theta_0 + \sum_{i=1}^N \lambda_i \cdot f_{\tau_i} \cdot \tau_i}{Z}$$

# Estimating Fisher Information

$$F_{\theta} = \mathbb{E}_{p_{\theta}(y|x)} \nabla_{\theta} \log p_{\theta}(y | x) \nabla_{\theta} \log p_{\theta}(y | x)^{\top}$$

Empirical and diagonal approximation:

$$f_{\theta} = \frac{1}{|\mathcal{D}|} \sum_{\mathcal{D}} (\nabla_{\theta} \log p(y | x))^2$$



# The algorithm

Step 1: Create task vectors

Step 2: Create Fisher vectors

Step 3: Merge

**Algorithm 1** Pseudocode for **removing hallucinations** and **promoting abstraction** with EWR. Note that we apply  $(\cdot)^2$  element-wise.

**Input** Dialogues  $\mathcal{D}$ , **hallucinated anti-expert dataset**  $\mathcal{D}^{\text{AE}}$ , **abstractive expert dataset**  $\mathcal{D}^{\text{E}}$ , initial parameter set  $\theta$

**Output**  $\theta'$

$$\theta_0 \leftarrow \text{finetune}(\theta, \mathcal{D})$$

$$\theta_{\text{AE}} \leftarrow \text{finetune}(\theta_0, \mathcal{D}^{\text{AE}})$$

$$\tau_1 \leftarrow \theta_{\text{AE}} - \theta_0$$

$$\theta_{\text{E}} \leftarrow \text{finetune}(\theta_0, \mathcal{D}^{\text{E}})$$

$$\tau_2 \leftarrow \theta_{\text{E}} - \theta_0$$

$$\mathbf{f}_{\theta_0} \leftarrow \frac{1}{|\mathcal{D}|} \sum_{\mathcal{D}} (\nabla \log p_{\theta_0}(u_{T+1} | u_1^T, \hat{\mathcal{K}}))^2$$

$$\mathbf{f}_{\tau_1} \leftarrow \frac{1}{|\mathcal{D}^{\text{AE}}|} \sum_{\mathcal{D}^{\text{AE}}} (\nabla \log p_{\tau_1}(u_{T+1} | u_1^T, \hat{\mathcal{K}}))^2$$

$$\mathbf{f}_{\tau_2} \leftarrow \frac{1}{|\mathcal{D}^{\text{E}}|} \sum_{\mathcal{D}^{\text{E}}} (\nabla \log p_{\tau_2}(u_{T+1} | u_1^T, \hat{\mathcal{K}}))^2$$

$$\theta' \leftarrow \frac{\lambda_0 \cdot \mathbf{f}_{\theta_0} \cdot \theta_0 - \lambda_1 \cdot \mathbf{f}_{\tau_1} \cdot \tau_1 + \lambda_2 \cdot \mathbf{f}_{\tau_2} \cdot \tau_2}{Z}$$



# Results

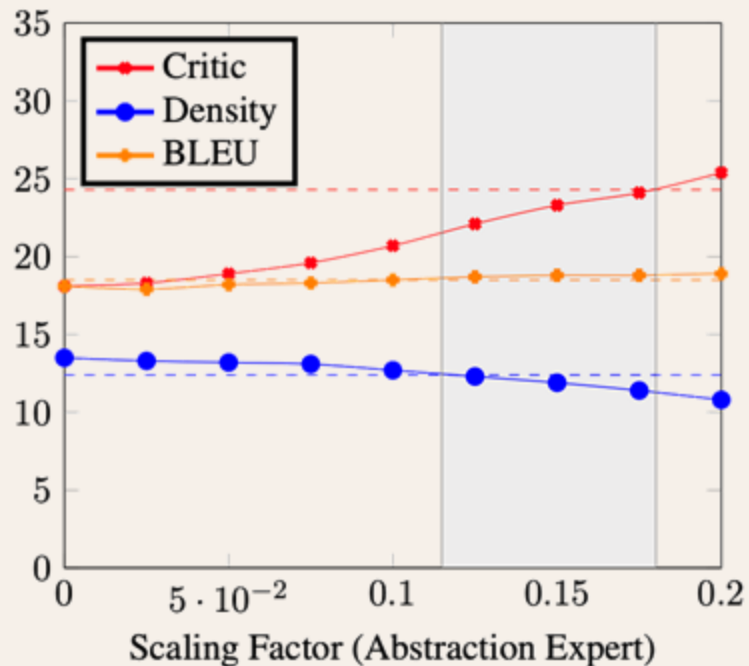
Model	BLEU( $\uparrow$ )	Critic( $\downarrow$ )	$Q^2$ ( $\uparrow$ )	BFI( $\uparrow$ )	F1( $\uparrow$ )
	$(y, \hat{y})$		$(y, \hat{\mathcal{K}})$		
	WoW <sub>unseen</sub>				
Flan-T5 <sub>base</sub>	18.1	22.7	74.0	84.8	78.7
+ TA	18.8	19.2	75.7	82.8	75.0
+ EWR	17.4 ( $\downarrow$ 0.7)	17.7 ( $\downarrow$ 5.0)	78.4 ( $\uparrow$ 4.4)	86.9 ( $\uparrow$ 2.1)	81.6 ( $\uparrow$ 2.9)
	DSTC11				
Flan-T5 <sub>base</sub>	7.9	76.6	49.7	54.6	37.1
+ TA	8.0	60.0	51.0	59.9	<b>43.6</b>
+ EWR	<b>9.6</b> ( $\uparrow$ 1.7)	<b>41.1</b> ( $\downarrow$ 35.5)	<b>57.3</b> ( $\uparrow$ 7.6)	<b>60.0</b> ( $\uparrow$ 5.4)	38.6 ( $\uparrow$ 1.5)
	FaithDial				
Flan-T5 <sub>base</sub>	15.1	0.3	<b>66.4</b>	80.9	73.7
+ TA	<b>15.3</b>	<b>0.1</b>	57.5	77.3	67.6
+ EWR	14.9 ( $\downarrow$ 0.2)	<b>0.1</b> ( $\downarrow$ 0.2)	<b>66.4</b> (-0.0)	<b>81.7</b> ( $\uparrow$ 0.8)	<b>75.0</b> ( $\uparrow$ 1.3)

Table 2: EWR improves faithfulness on unseen topics (WoW<sub>unseen</sub>), multi-document corpora (DSTC11), and datasets with cleaned ground-truth annotations (Faith-

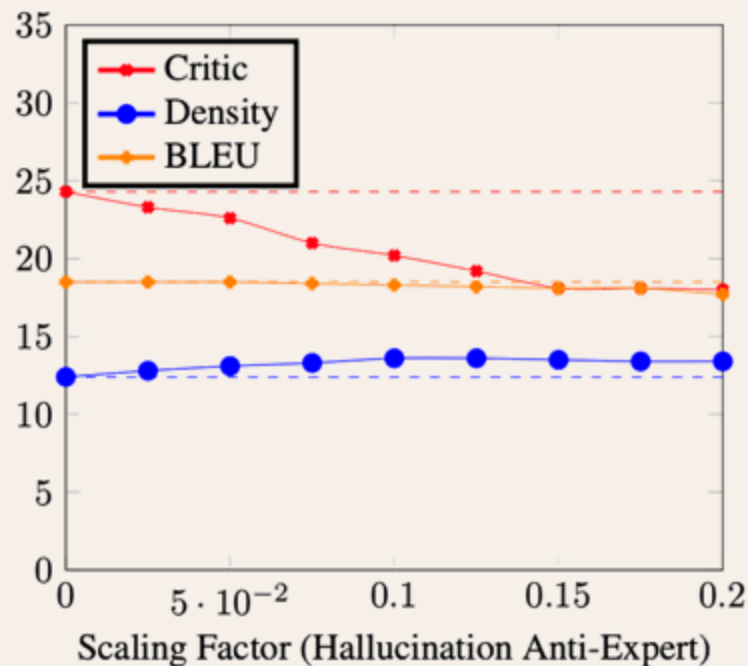
Model	WoW			DSTC9		
	A ( $\uparrow$ )	C ( $\uparrow$ )	P ( $\uparrow$ )	A ( $\uparrow$ )	C ( $\uparrow$ )	P ( $\uparrow$ )
Flan-T5 <sub>base</sub>	72.3	1.74	1.19	89.7	2.83	1.71
+ EWR <sub>abs</sub>	<b>75.1</b>	1.62	1.25	<b>94.7*</b>	2.41	1.49
CTRL	85.5*	1.58	1.12	94.7*	2.72	1.42
+ TA	88.8*	1.58	1.16	97.0*	2.63	1.40
+ EWR	96.8 $\dagger$	1.50	1.08	98.0 $\dagger$	2.50	1.36
Quark	93.1 $\dagger$	1.51	1.05	86.0	2.89	1.66

Table 3: Human evaluation on 218 examples annotated by 3 expert annotators each. We measure attributability (A), Co-cooperativeness (C), and paraphrasing (P). \* indicates significant improvements wrt. Flan-T5<sub>base</sub> and  $\dagger$  wrt. to the next best method with  $p < 0.05$ .

# The Sweet Spot



(a) Faithfulness-Abtractiveness Trade-Off



(b) Faithfulness-Performance Trade-Off

# Code

<https://github.com/ndaheim/faithful-dialogue>



Baseline	<code>document_grounded_generation</code>	(context, knowledge) -> response
CTRL	<code>document_grounded_generation_ctrl</code>	(context, knowledge, control tokens) -> response
Quark	<code>document_grounded_generation_Quark</code>	CTRL tokens as quantized reward for samples drawn from the model during training
Noisy Channel Reranking	<code>noisy_channel_reranking</code>	reranking controllably with a faithfulness and fluency expert
DExperts	<code>document_grounded_generation_density_ratio</code>	faithfulness expert and anti-expert combined at inference time
Task Arithmetic	<code>document_grounded_generation</code>	task vector subtracted from base model using hallucination anti-expert
CaPE	<code>document_grounded_generation</code>	task vector subtracted from base model using hallucination anti-expert and faithfulness expert
EWR	<code>document_grounded_generation</code>	task vector subtracted from base model using hallucination anti-expert weighted by Fisher Information

Note that TA, CaPE and EWR just change the model parameters via interpolation, but not the model architecture!

## Datasets

1. Wizard-of-Wikipedia
2. FaithDial
3. DSTC9
4. DSTC11

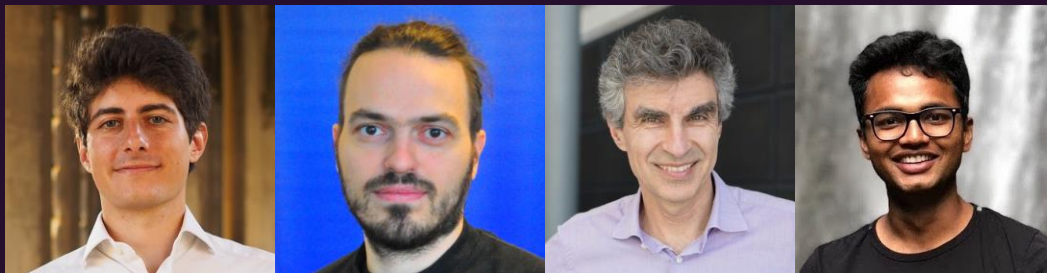
## Metrics

1. BLEU
2. BERTScore
3. Faithfulness Critic
4.  $Q^2$
5. Knowledge F1
6. Density & Coverage

# *Polytropon*: Joint Routing and Adaptation

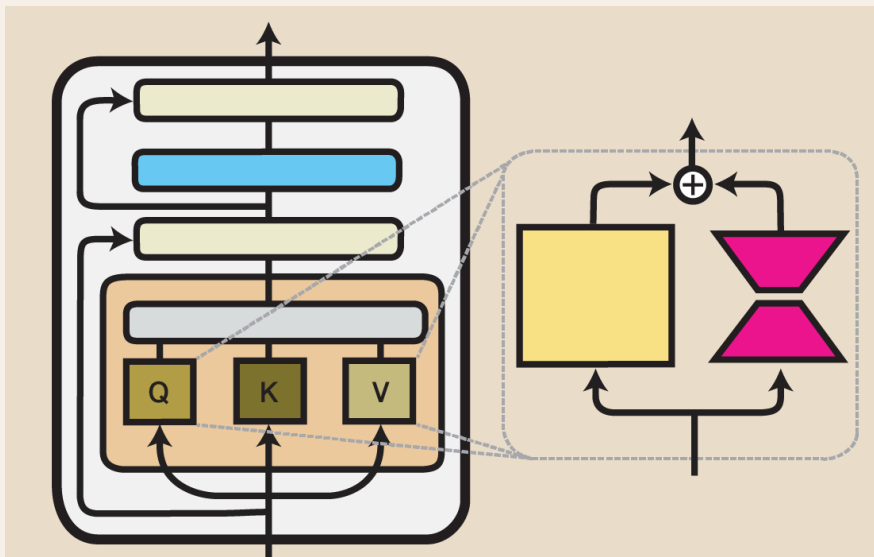
Edoardo M. Ponti, Alessandro Sordoni, Yoshua Bengio, Siva Reddy.  
*Combining Parameter-efficient Modules for Task-level Generalisation.*

EACL 23



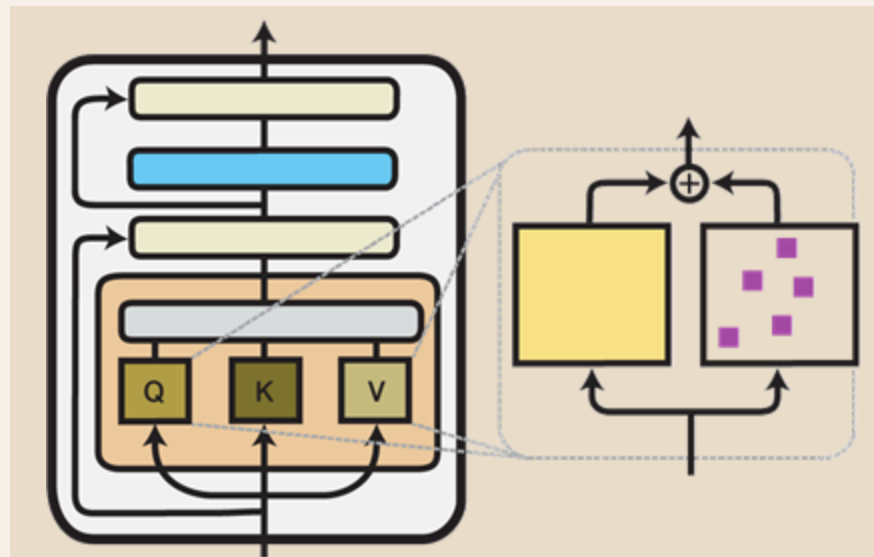
**Goal: adapting general-purpose LLMs  
efficiently and systematically to new tasks**

# An inventory of Modules (=Adapters)



**Low-rank Adapter** [Hu et al. 2021]

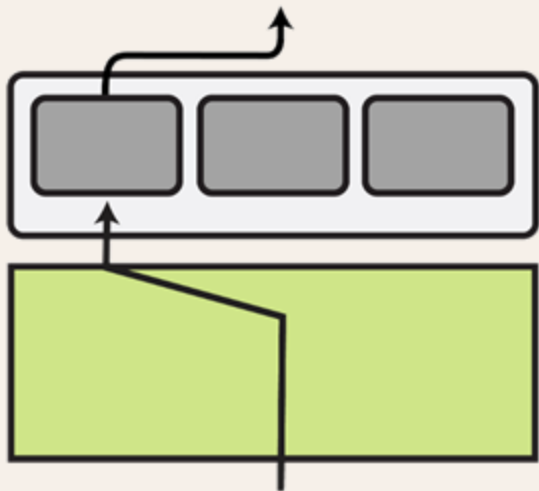
$$\phi_j = W + A_j^\top B_j$$
$$A, B \in \mathbb{R}^{d \times r} \quad r \ll d$$



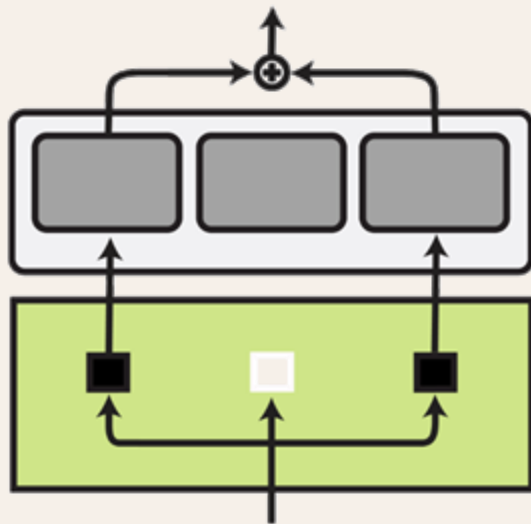
**Sparse Adapter** [Ansell et al. 2021]

$$\phi_j = W + A_j$$
$$A \in \mathbb{R}^{d \times d} \text{ is sparse}$$

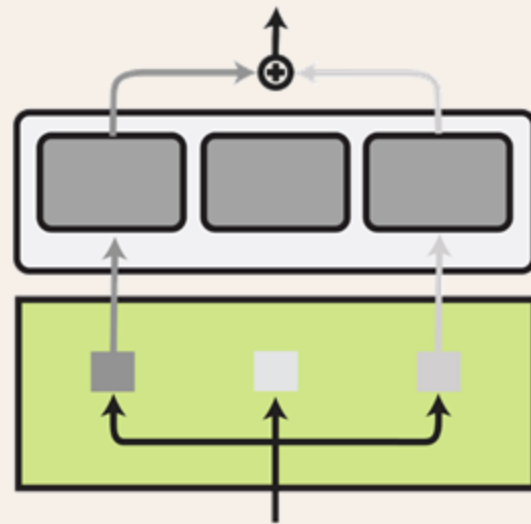
# Routing: Learned and Variable-size



(a) Fixed Routing



(b) Learned Routing (Hard)



(c) Learned Routing (Soft)

$$\alpha_{i,j} = \text{sigmoid} \left[ \log \frac{\text{sigmoid}(\hat{\alpha}_{i,j}) u}{(1 - \text{sigmoid}(\hat{\alpha}_{i,j})) (1 - u)} \right]^{1/\tau}$$

$$u \sim \text{Uniform}(0, 1).$$

# Polytropon: Discovering Skills End-to-end

Core idea: jointly learn *adapters* (modules) and *variable-size routing* to fine-tune a LLM.

$$f'_i(x) = \sum_j \frac{\alpha_j}{Z} f(x; \theta, \phi_j)$$

modules  $\Theta$

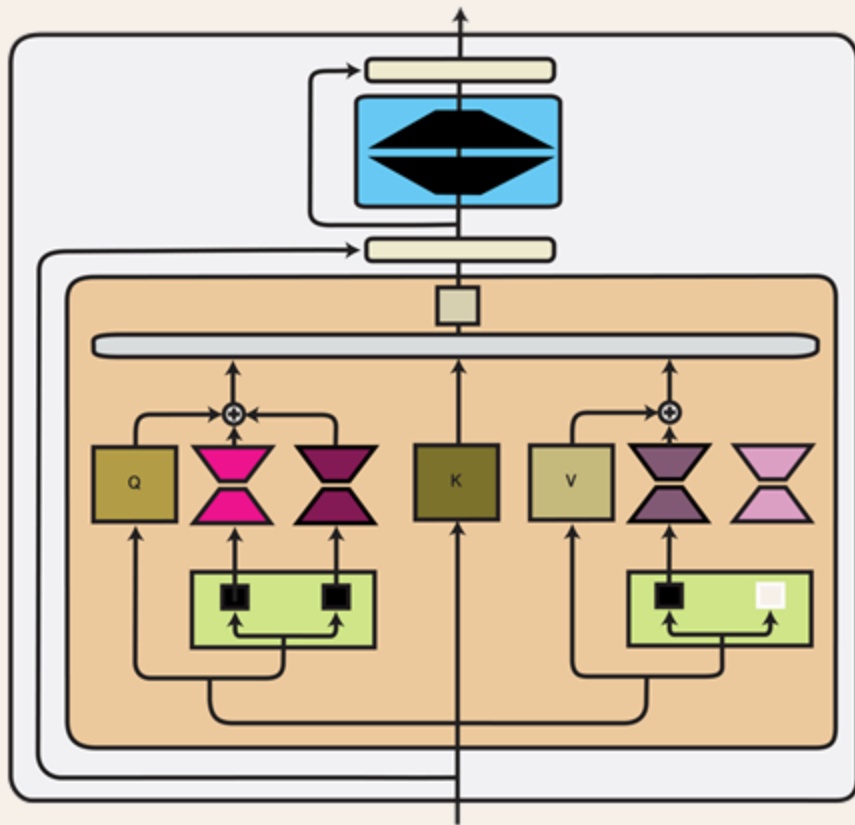
skills

Adapter 1
Adapter 2
Adapter 3

routing  $Z$

tasks


skills





# Comparison with Mixture-of-Experts

	<b>Polytropon</b>	<b>MoEs</b>
<b>Purpose</b>	Few-shot generalisation	Scaling LLMs with sparsity
<b>Training</b>	Fine-tuning	Pre-training
<b>Modules</b>	Adapters	FFNs
<b>Routing</b>	Variable-size, task-level	Top- $k$ , token-level

# Challenges of Learned Routing

- **Training Instability**

Router and modules are untrained → routing dynamics never stabilise.

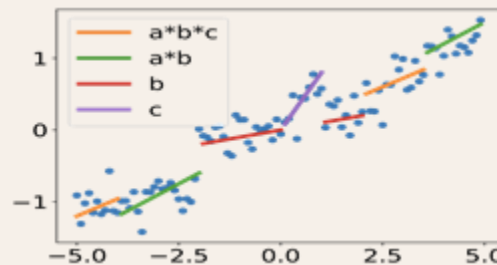
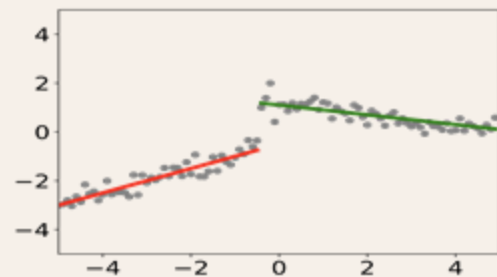
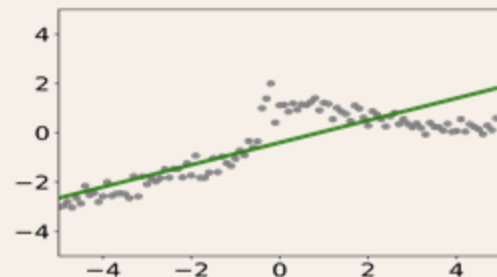
- **Underfitting (aka module collapse)**

The router falls into a local optimum, choosing a few modules exclusively

- **Overfitting**

Risk of overfitting to the noise.

Plots courtesy of [Rosenbaum et al. \(2017\)](#)



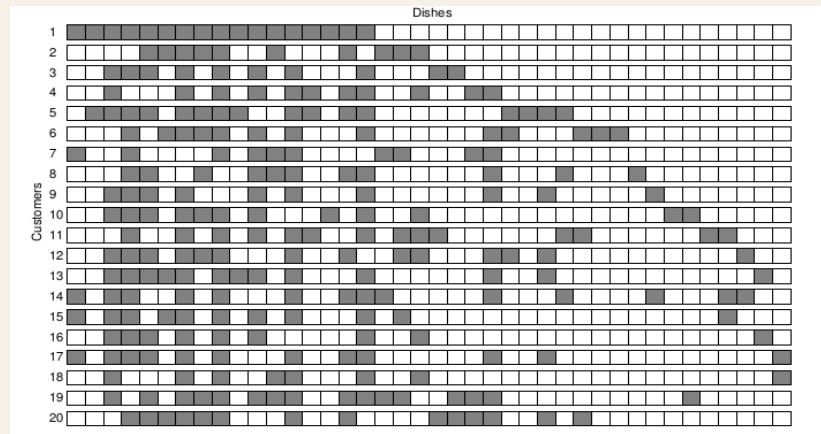
# Inductive Biases

## Module collapse [Rosenbaum et al. 2019]

Only a small number of modules from the inventory are selected.  
Results from excessively favouring *exploitation over exploration*.

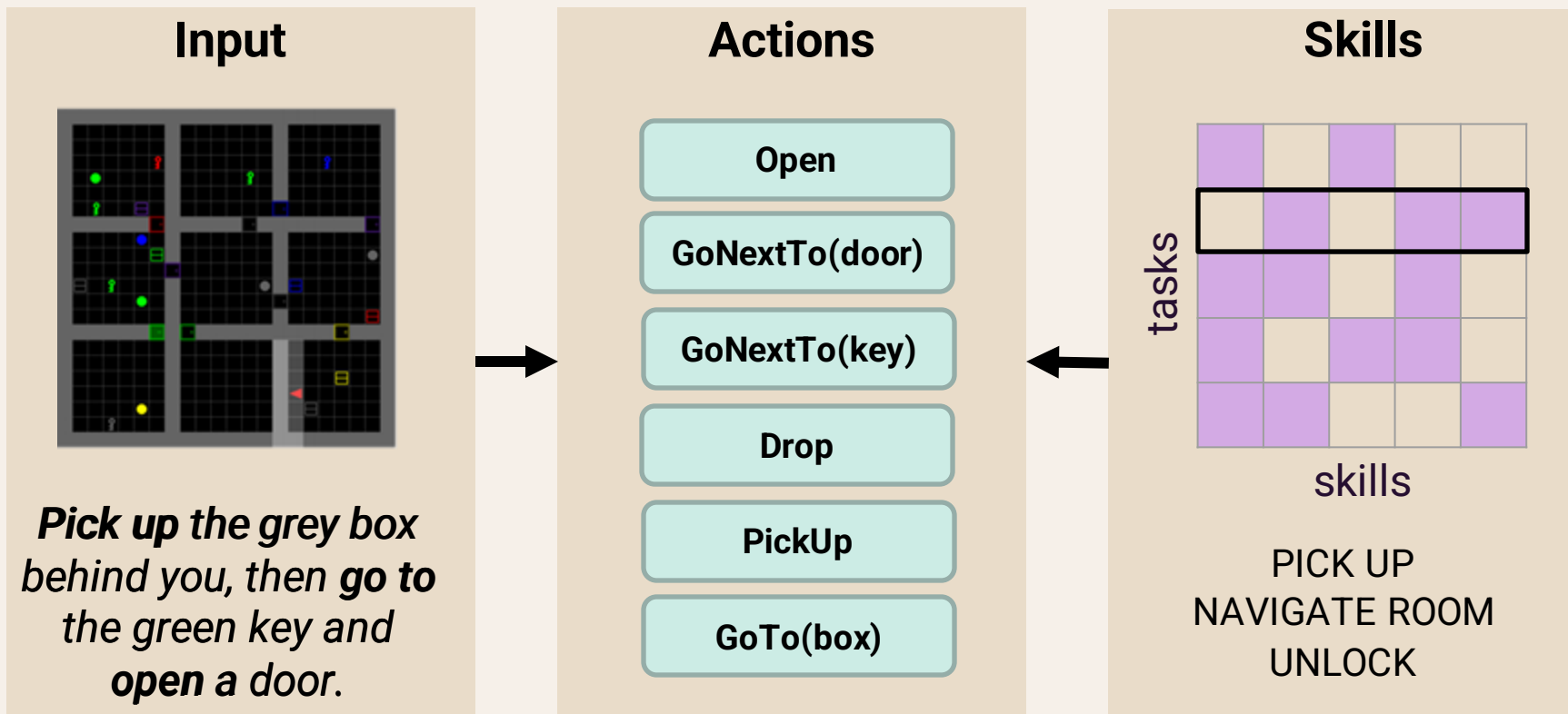
## Indian Buffet Process [Griffiths and Ghahramani 2011]

## Dual-speed Learning Rate

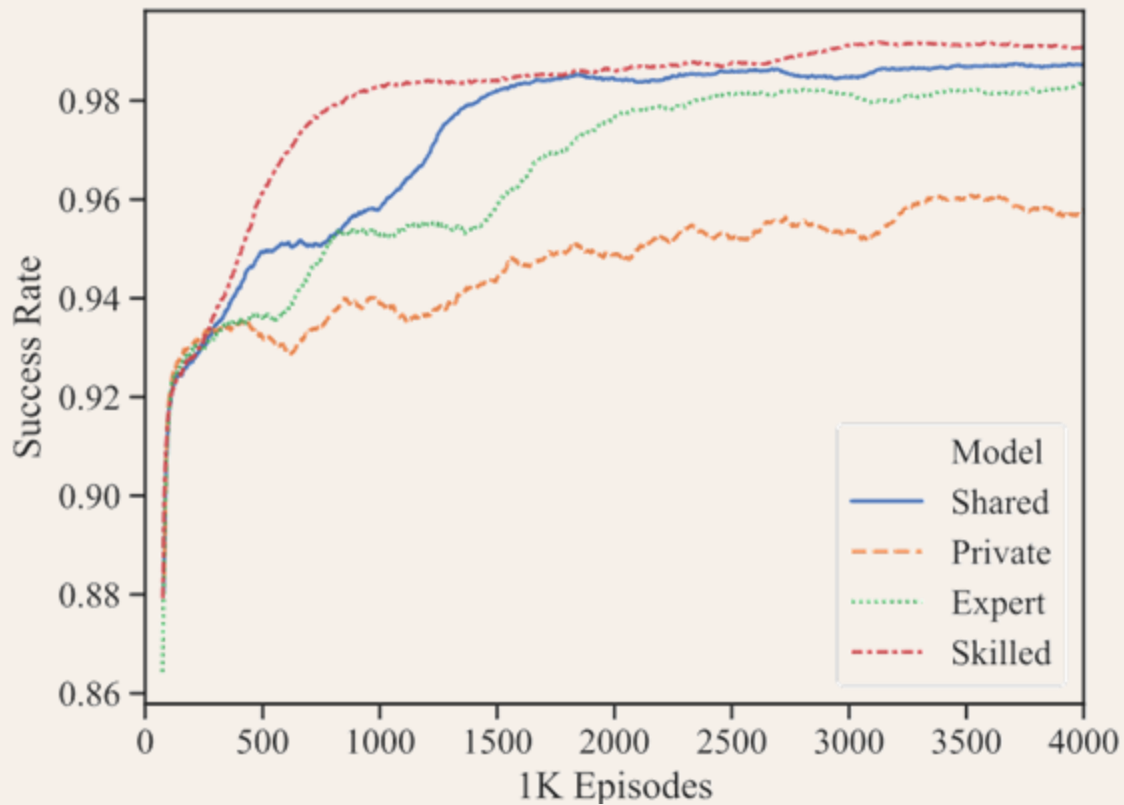


Higher for the routing function  
than the module parameters.  
Intuition: *coarse-to-fine dynamic*.

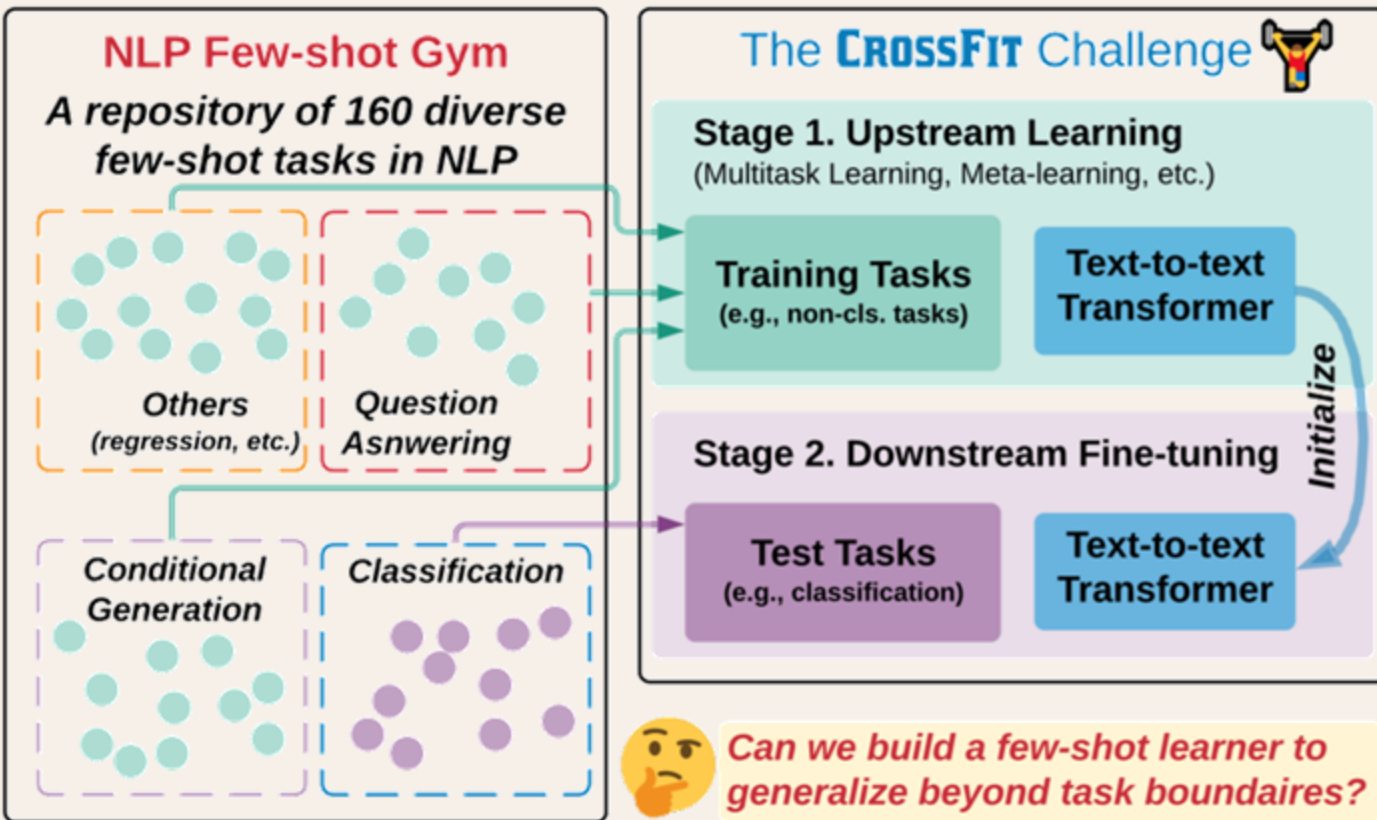
# Instruction Following in RL: BabyAI [Chevalier-Boisvert et al. 2018]



# Reinforcement Learning Results: Sample Efficiency

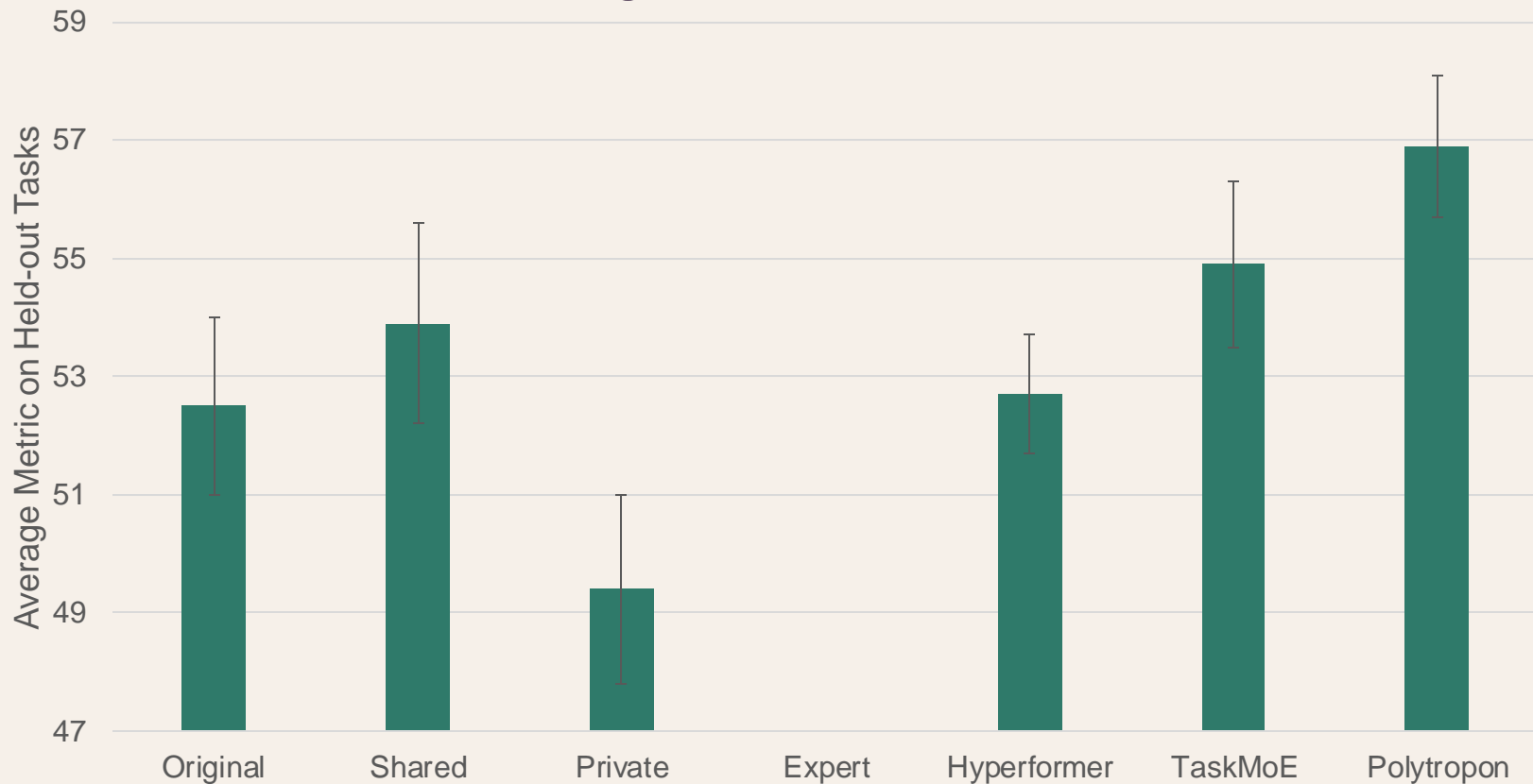


# NLP: CrossFit [Ye et al. 2021]

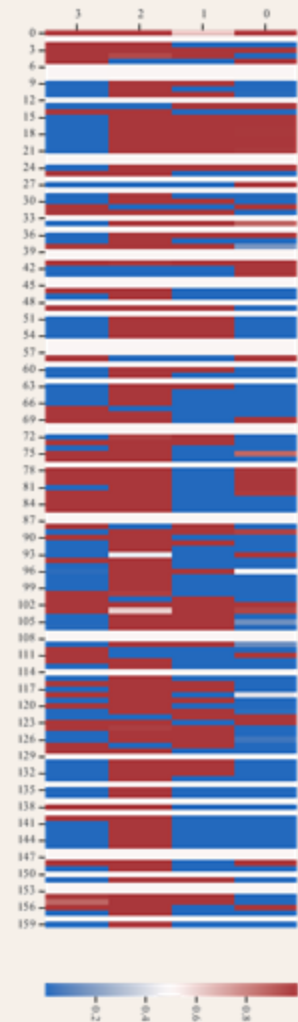
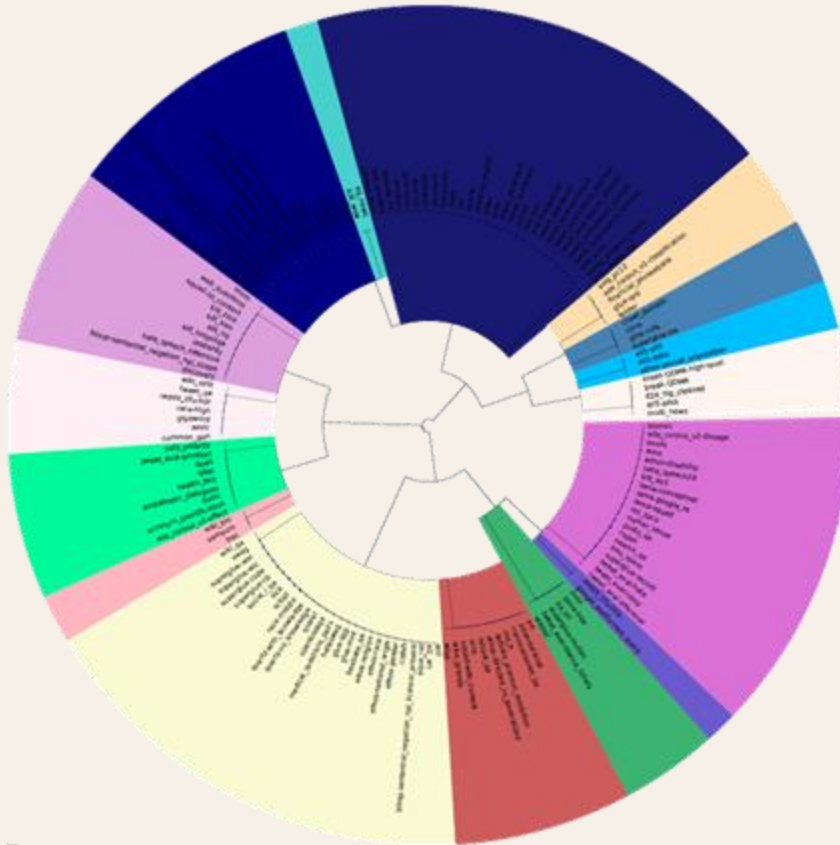


# NLP Results:

## Few-shot Learning in Unseen Tasks



# Interpretability: Task Hierarchy





# Conclusions

- Efficient multi-task learning by **implementing *skills* through adapters.**
- Inductive biases that encourage module (re)combination, e.g. allowing for **variable-size** module **routing**
- **Higher sample efficiency** in multi-task reinforcement learning and **better few-shot adaptation** in multi-task supervised learning

# Code

<https://github.com/microsoft/mttl>



## MTTL

---

MTTL - Multi-Task Transfer Learning

### Setup

---

Install Python packages:

```
pip install -r requirements.txt
```

The package `promptsource` currently requires Python 3.7. Alternative versions require local installations (see their [documentation](#)).

Download the datasets:

```
bash scripts/create_datasets.sh
```

### Multi-task Pre-training

---

The general command:

```
python pl_train.py -c $CONFIG_FILES -k $KWARGS
```

Multiple `CONFIG_FILES` can be concatenated as `file1+file2`. To modify defaults, `KWARGS` can be expressed as `key=value`.

### Test Fine-Tuning

---

To perform finetuning for a test task, use the script `pl_finetune.py`

### Hyper-parameter Search for Test Fine-Tuning

---

To perform an hyperparameter search for a test task, use the script `pl_finetune_tune.py`. The script will just call the functions in `pl_finetune.py` in a loop. The script itself defines hp ranges for different fine-tuning types.

Thanks for your attention!