



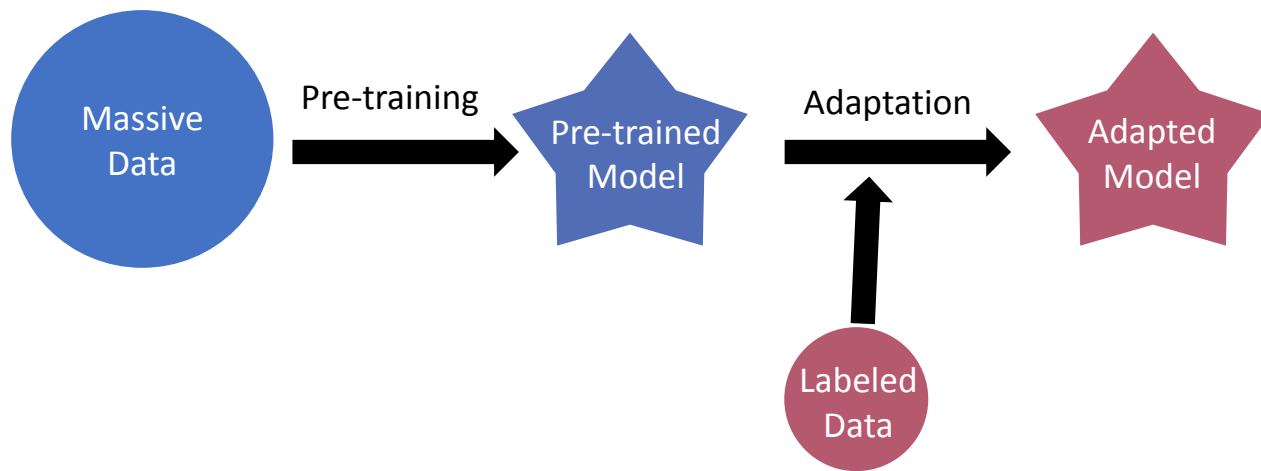
Towards Few-Shot Adaptation of Foundation Models via Multitask Finetuning

Zhuoyan Xu, Zhenmei Shi, Junyi Wei, Fangzhou Mu, Yin Li, Yingyu Liang
University of Wisconsin - Madison



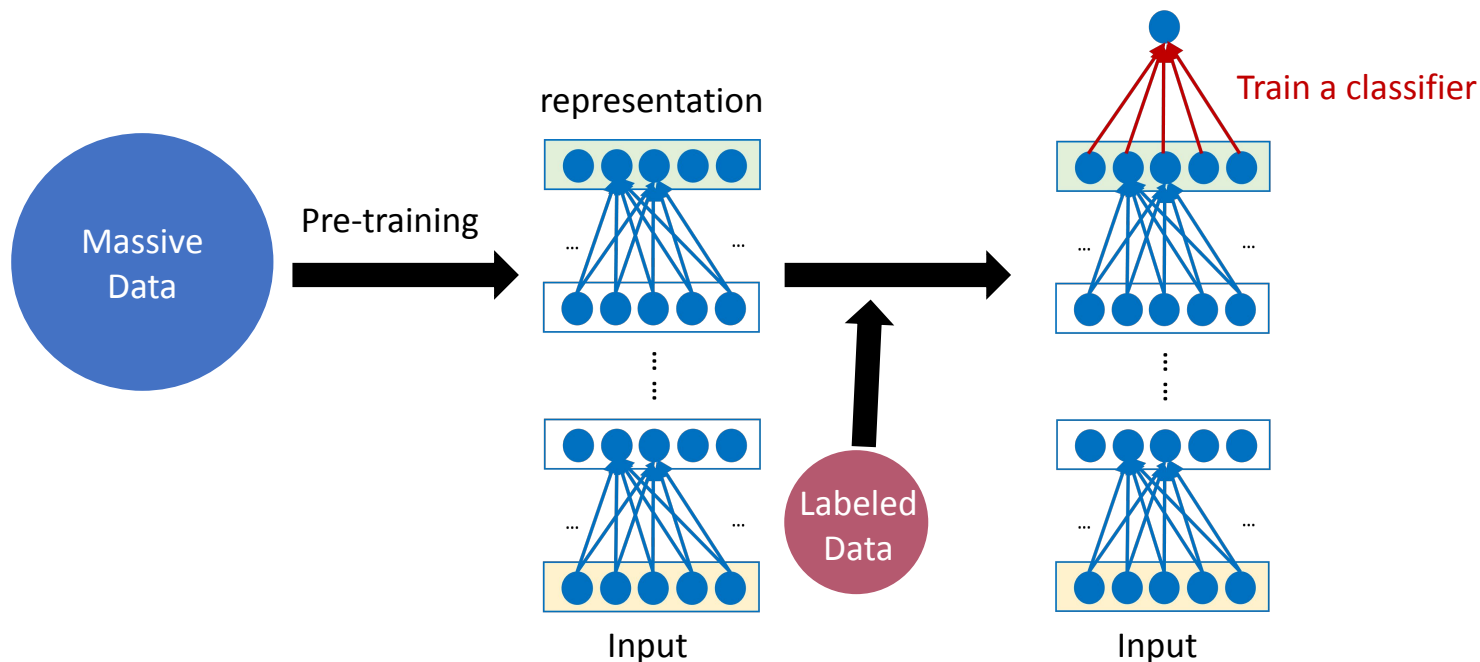
New Paradigm: Pre-trained Representations

Paradigm shift: supervised learning \Rightarrow pre-training + adaptation



New Paradigm: Pre-trained Representations

Paradigm shift: supervised learning \longrightarrow pre-training + adaptation



New Paradigm: Pre-trained Representations

Paradigm shift: supervised learning \longrightarrow pre-training + adaptation

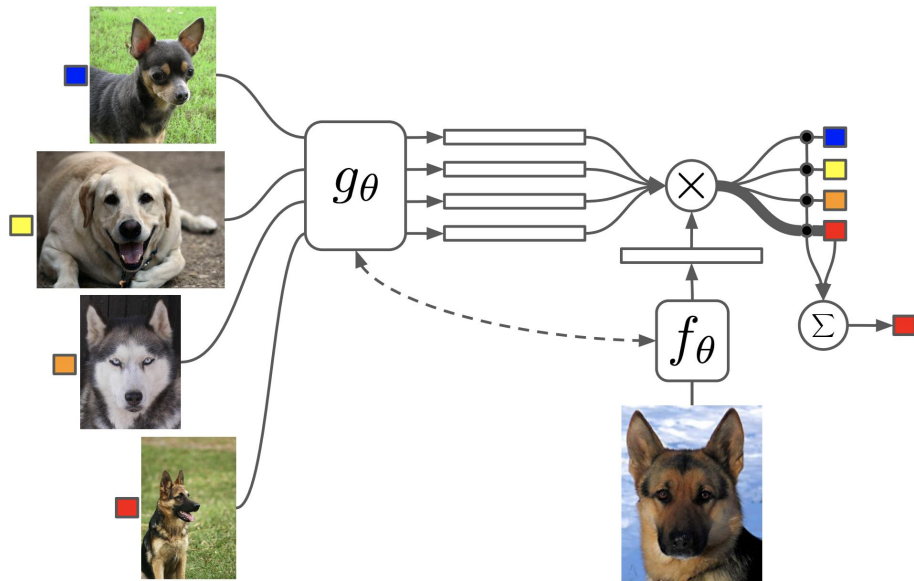
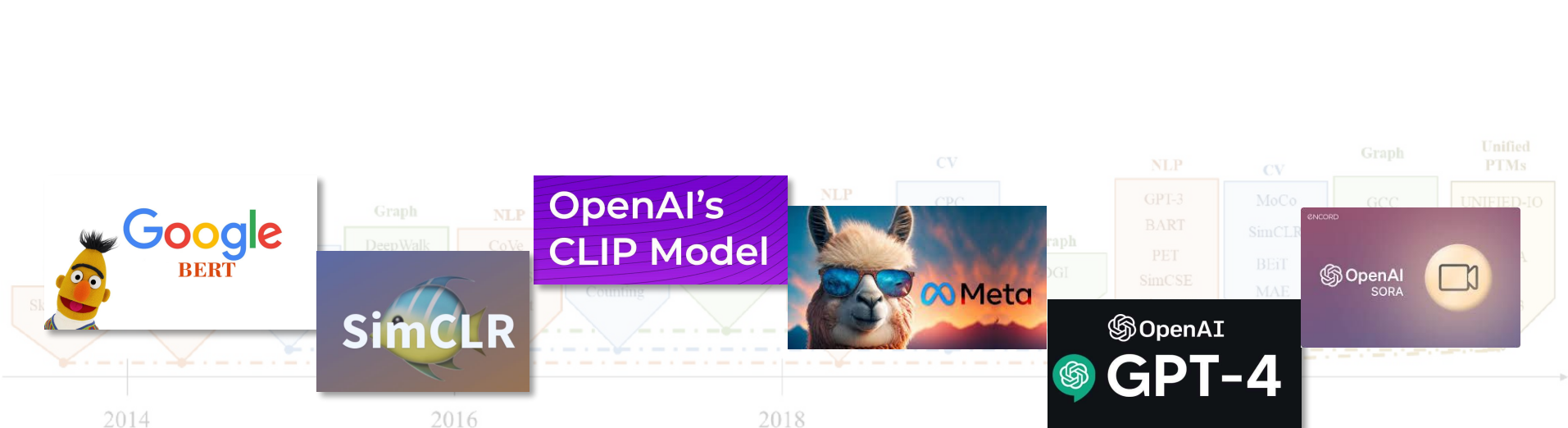


Figure 1: Matching Networks architecture

Adaptation of a pre-trained image encoder

Figures from: *Matching Networks for One Shot Learning*, 2017.

Intro - Foundation Model

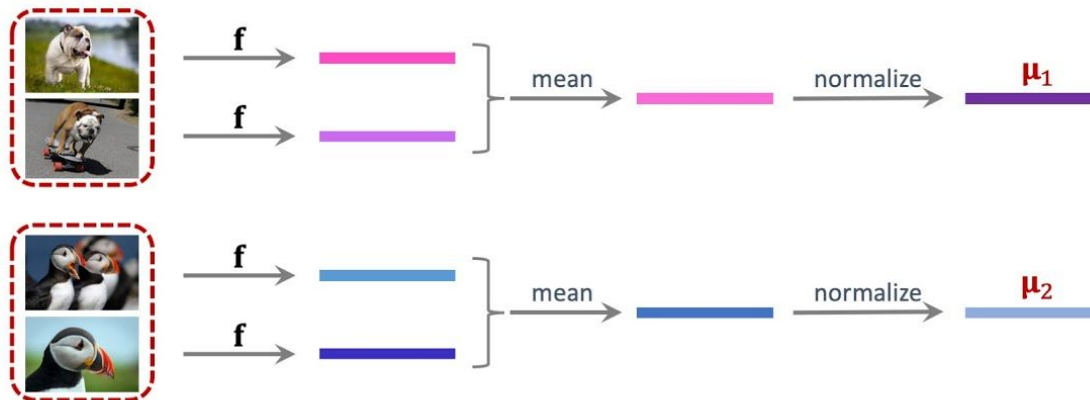


The history and evolution of foundation models

Figures from: *A Comprehensive Survey on Pretrained Foundation Models: A History from BERT to ChatGPT, 2023.*

Intro - Foundation Model

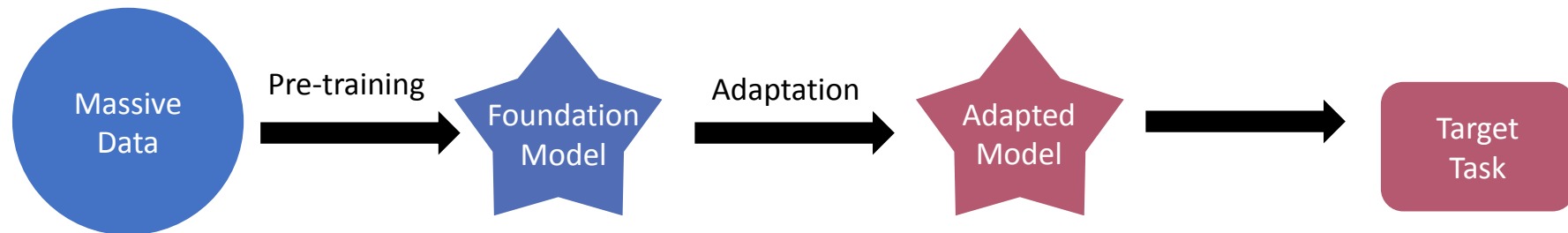
Few-Shot Learning: Pretraining + Fine Tuning



Label Efficiency

Figures from: https://www.youtube.com/watch?v=U6uFOIURcd0&ab_channel=ShusenWang, 2020

Paradigm: Pre-training + Adaptation



Pre-training

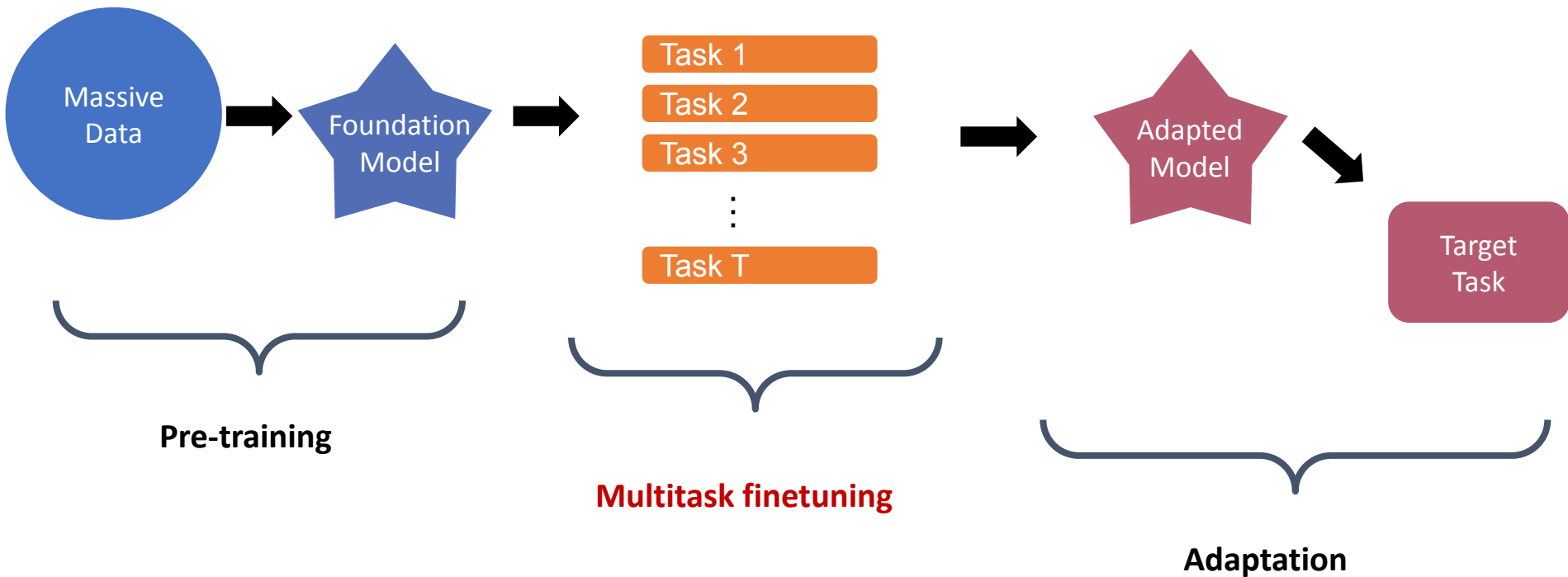


Adaptation



Q: Can we improve this?

Pre-training + Finetuning + Adaptation



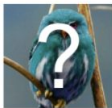
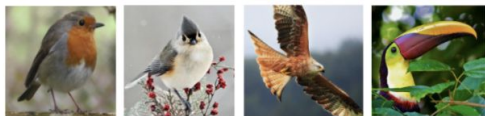
Training

Train dataset #1: "cat-bird"

cats



birds



Train dataset #2: "flower-bike"

flowers



bikes



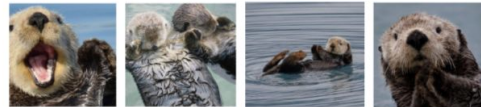
Testing

Test dataset: "dog-otter"

dogs



otters



An example of 4-shot 2-class image classification

Figures from: [Meta-Learning: Learning to Learn Fast](#), 2018.

Diversity and Consistency

Definition 1 (**Diversity and Consistency (Informal)**)

Consider the latent feature space of target task data and finetuning task data.

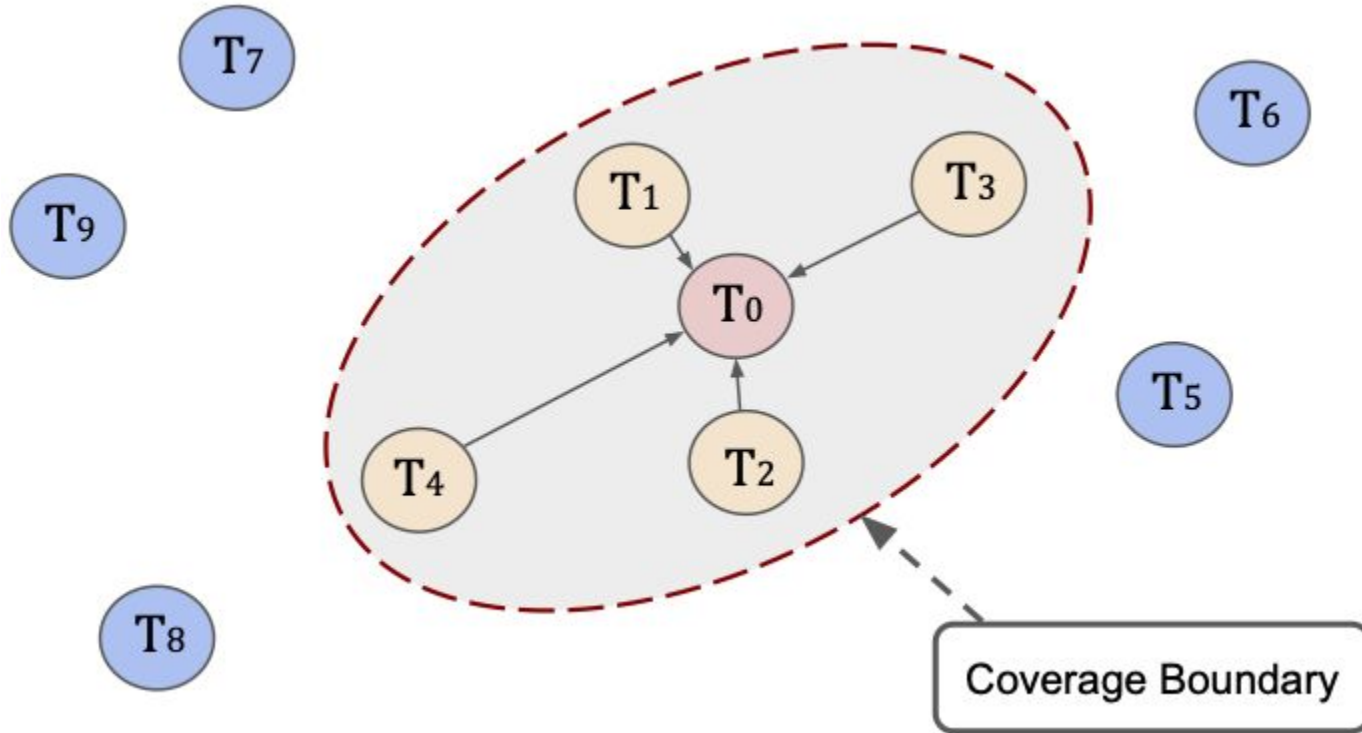
Diversity \rightarrow coverage.

Consistency \rightarrow similarity.

Theorem (**Multitask finetuning loss (Informal)**)

Given pre-trained foundation model and certain loss on target task, employing multitask finetuning on the this pretrained model with sufficient tasks and samples can help further reduce the loss on the target task.

Practical solution: Task selection



Experiments: Task selection algorithm

Pretrained	Selection	INet	Omglot	Acrafft	CUB	QDraw	Fungi	Flower	Sign	COCO
CLIP	Random	56.29	65.45	31.31	59.22	36.74	31.03	75.17	33.21	30.16
	No Con.	60.89	72.18	31.50	66.73	40.68	35.17	81.03	37.67	34.28
	No Div.	56.85	73.02	32.53	65.33	40.99	33.10	80.54	34.76	31.24
	Selected	60.89	74.33	33.12	69.07	41.44	36.71	80.28	38.08	34.52
DINOv2	Random	83.05	62.05	36.75	93.75	39.40	52.68	98.57	31.54	47.35
	No Con.	83.21	76.05	36.32	93.96	50.76	53.01	98.58	34.22	47.11
	No Div.	82.82	79.23	36.33	93.96	55.18	52.98	98.59	35.67	44.89
	Selected	83.21	81.74	37.01	94.10	55.39	53.37	98.65	36.46	48.08
MoCo v3	Random	59.66	60.72	18.57	39.80	40.39	32.79	58.42	33.38	32.98
	No Con.	59.80	60.79	18.75	40.41	40.98	32.80	59.55	34.01	33.41
	No Div.	59.57	63.00	18.65	40.36	41.04	32.80	58.67	34.03	33.67
	Selected	59.80	63.17	18.80	40.74	41.49	33.02	59.64	34.31	33.86

Table 1: Results evaluating our task selection algorithm on Meta-dataset using ViT-B backbone. No Con.: Ignore consistency. No Div.: Ignore diversity. Random: Ignore both consistency and diversity.

Experiments: Effectiveness of Multitask Finetuning

pretrained	backbone	method	miniImageNet		tieredImageNet		DomainNet	
			1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
MoCo v3	ViT-B	Adaptation	75.33 (0.30)	92.78 (0.10)	62.17 (0.36)	83.42 (0.23)	24.84 (0.25)	44.32 (0.29)
		Standard FT	75.38 (0.30)	92.80 (0.10)	62.28 (0.36)	83.49 (0.23)	25.10 (0.25)	44.76 (0.27)
		Ours	80.62 (0.26)	93.89 (0.09)	68.32 (0.35)	85.49 (0.22)	32.88 (0.29)	54.17 (0.30)
	ResNet50	Adaptation	68.80 (0.30)	88.23 (0.13)	55.15 (0.34)	76.00 (0.26)	27.34 (0.27)	47.50 (0.28)
		Standard FT	68.85 (0.30)	88.23 (0.13)	55.23 (0.34)	76.07 (0.26)	27.43 (0.27)	47.65 (0.28)
		Ours	71.16 (0.29)	89.31 (0.12)	58.51 (0.35)	78.41 (0.25)	33.53 (0.30)	55.82 (0.29)
DINO v2	ViT-S	Adaptation	85.90 (0.22)	95.58 (0.08)	74.54 (0.32)	89.20 (0.19)	52.28 (0.39)	72.98 (0.28)
		Standard FT	86.75 (0.22)	95.76 (0.08)	74.84 (0.32)	89.30 (0.19)	54.48 (0.39)	74.50 (0.28)
		Ours	88.70 (0.22)	96.08 (0.08)	77.78 (0.32)	90.23 (0.18)	61.57 (0.40)	77.97 (0.27)
	ViT-B	Adaptation	90.61 (0.19)	97.20 (0.06)	82.33 (0.30)	92.90 (0.16)	61.65 (0.41)	79.34 (0.25)
		Standard FT	91.07 (0.19)	97.32 (0.06)	82.40 (0.30)	93.07 (0.16)	61.84 (0.39)	79.63 (0.25)
		Ours	92.77 (0.18)	97.68 (0.06)	84.74 (0.30)	93.65 (0.16)	68.22 (0.40)	82.62 (0.24)
Supervised pretraining on ImageNet	ViT-B	Adaptation	94.06 (0.15)	97.88 (0.05)	83.82 (0.29)	93.65 (0.13)	28.70 (0.29)	49.70 (0.28)
		Standard FT	95.28 (0.13)	98.33 (0.04)	86.44 (0.27)	94.91 (0.12)	30.93 (0.31)	52.14 (0.29)
		Ours	96.91 (0.11)	98.76 (0.04)	89.97 (0.25)	95.84 (0.11)	48.02 (0.38)	67.25 (0.29)
	ResNet50	Adaptation	81.74 (0.24)	94.08 (0.09)	65.98 (0.34)	84.14 (0.21)	27.32 (0.27)	46.67 (0.28)
		Standard FT	84.10 (0.22)	94.81 (0.09)	74.48 (0.33)	88.35 (0.19)	34.10 (0.31)	55.08 (0.29)
		Ours	87.61 (0.20)	95.92 (0.07)	77.74 (0.32)	89.77 (0.17)	39.09 (0.34)	60.60 (0.29)

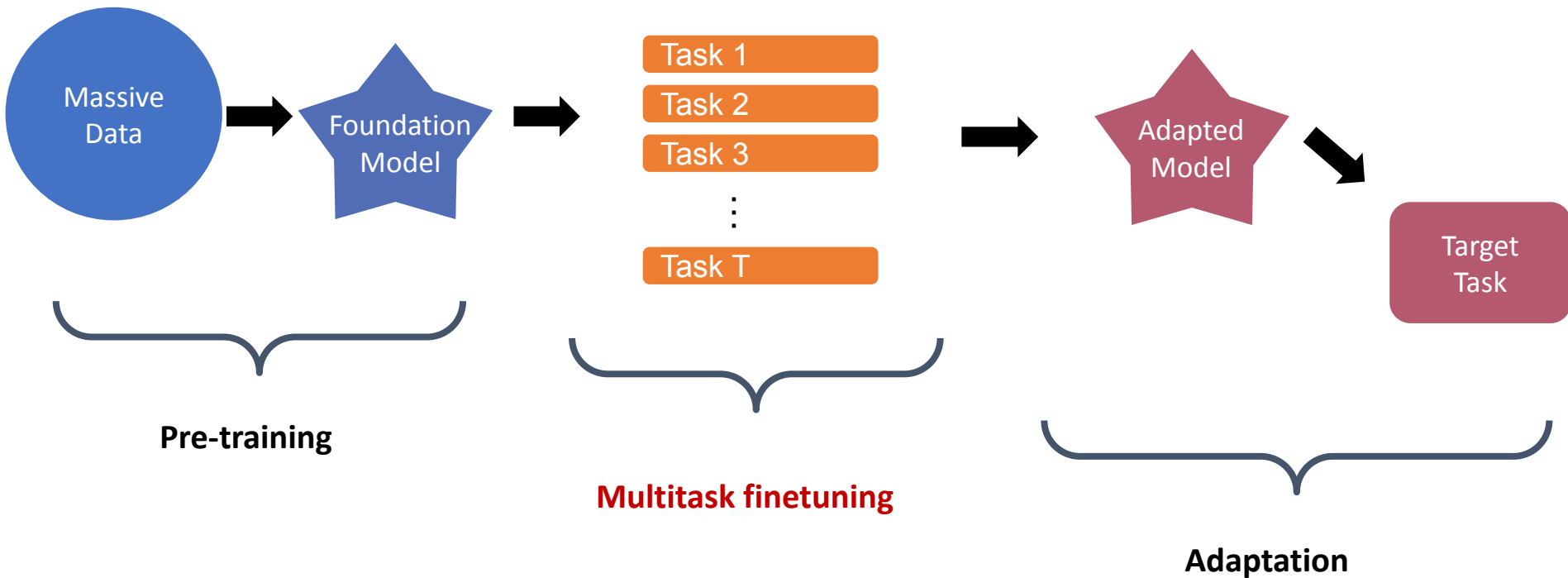
Table 2: **Results of few-shot image classification.** We report average classification accuracy (%) with 95% confidence intervals on test splits. Adaptation: Direction adaptation without finetuning; Standard FT: Standard finetuning; Ours: Our multitask finetuning; 1-/5-shot: number of labeled images per class in the target task.

Experiments: Few-shot Language task

	SST-2 (acc)	SST-5 (acc)	MR (acc)	CR (acc)	MPQA (acc)	Subj (acc)	TREC (acc)	CoLA (Matt.)
Prompt-based zero-shot	83.6	35.0	80.8	79.5	67.6	51.4	32.0	2.0
Multitask FT zero-shot	92.9	37.2	86.5	88.8	73.9	55.3	36.8	-0.065
+ task selection	92.5	34.2	87.1	88.7	71.8	72.0	36.8	0.001
Prompt-based FT [†]	92.7 (0.9)	47.4 (2.5)	87.0 (1.2)	90.3 (1.0)	84.7 (2.2)	91.2 (1.1)	84.8 (5.1)	9.3 (7.3)
Multitask Prompt-based FT	92.0 (1.2)	48.5 (1.2)	86.9 (2.2)	90.5 (1.3)	86.0 (1.6)	89.9 (2.9)	83.6 (4.4)	5.1 (3.8)
+ task selection	92.6 (0.5)	47.1 (2.3)	87.2 (1.6)	91.6 (0.9)	85.2 (1.0)	90.7 (1.6)	87.6 (3.5)	3.8 (3.2)
	MNLI (acc)	MNLI-mm (acc)	SNLI (acc)	QNLI (acc)	RTE (acc)	MRPC (F1)	QQP (F1)	
Prompt-based zero-shot	50.8	51.7	49.5	50.8	51.3	61.9	49.7	
Multitask FT zero-shot	63.2	65.7	61.8	65.8	74.0	81.6	63.4	
+ task selection	62.4	64.5	65.5	61.6	64.3	75.4	57.6	
Prompt-based FT [†]	68.3 (2.3)	70.5 (1.9)	77.2 (3.7)	64.5 (4.2)	69.1 (3.6)	74.5 (5.3)	65.5 (5.3)	
Multitask Prompt-based FT	70.9 (1.5)	73.4 (1.4)	78.7 (2.0)	71.7 (2.2)	74.0 (2.5)	79.5 (4.8)	67.9 (1.6)	
+ task selection	73.5 (1.6)	75.8 (1.5)	77.4 (1.6)	72.0 (1.6)	70.0 (1.6)	76.0 (6.8)	69.8 (1.7)	

Table 18: **Results of few-shot learning with NLP benchmarks.** All results are obtained using RoBERTa-large. We report the mean (and standard deviation) of metrics over 5 different splits. †: Result in [Gao et al. \(2021a\)](#) in our paper; FT: finetuning; task selection: select multitask data from customized datasets.

Take Home Message



Thanks!