

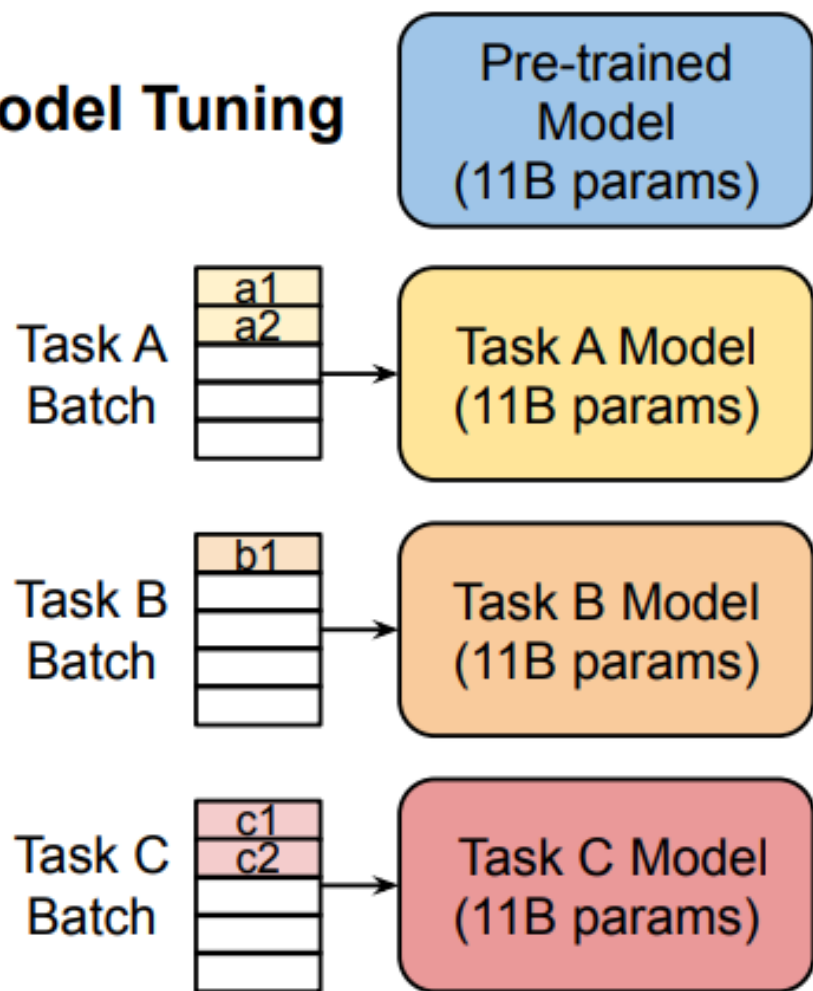
Prompt-tuning: сравнение методов

Гвоздева Дарья БПМИ211

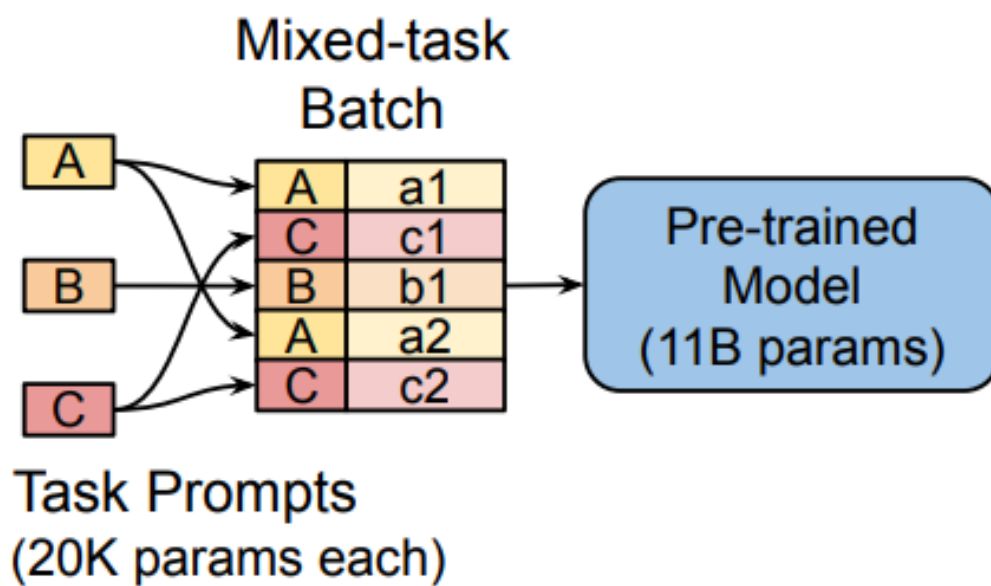
План на сегодня

- Идея prompt-tuning'a
- Parameter-Efficient Prompt Tuning
 - Суть подхода
 - Результаты эксперимента
 - Сравнение с прочими подходами
- Pre-trained Prompt Tuning
 - Суть подхода
 - Эффективность
- Итоги

Model Tuning



Prompt Tuning



Основная идея

$$\Pr(y|X) \rightarrow Pr_{\theta, \theta_p}(Y|[P; X])$$

Основная идея

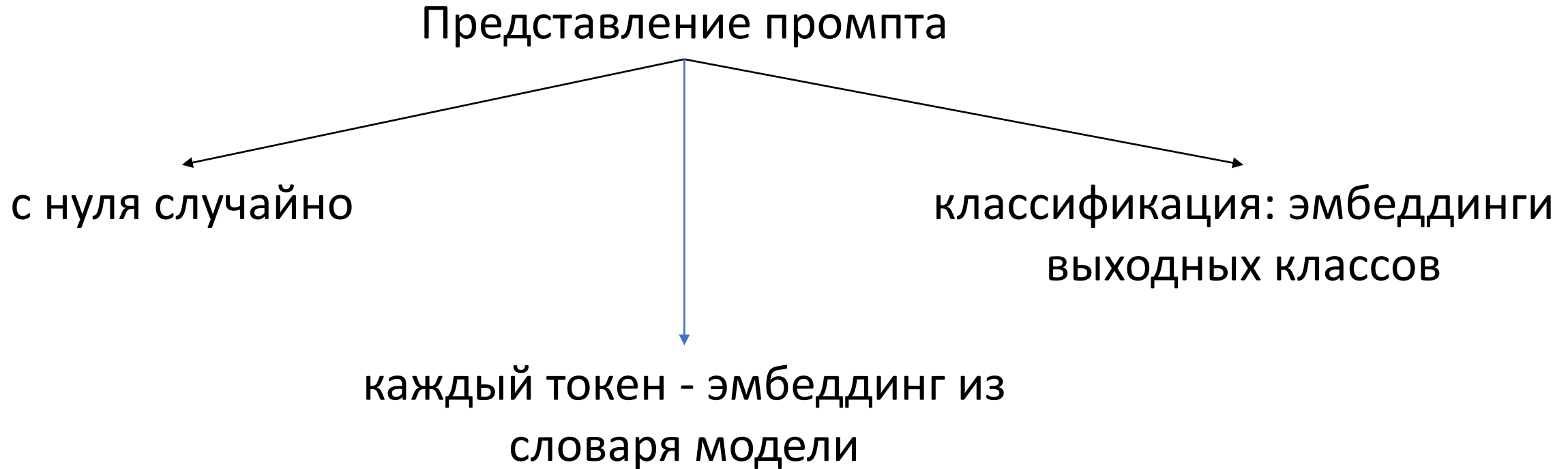
$$\Pr(y|X) \rightarrow Pr_{\theta, \theta_p}(Y|[P; X])$$

GPT3: $P = \{p_1, \dots, p_n\}$ - p_i меняются, θ fixed

PT: P (prompt)- *fixed*, θ_p fixed & own for each prompt, embeddings updated

Prompt design: обновляемые p_i из фиксированных эмбеддингов

Подходы к дизайну инициализации



Результаты

Setup:

- T5 различных размеров
- Дефолтная конфигурация в 100 токенов
- Benchmark: SuperGLUE
 - Каждый промпт – своя задача
 - Без мультитаска/смешивания данных
- 30.000 steps
- T5 -> стандартное отклонение + кросс энтропия
- lr: const 0.3
- Batchsize: 32
- Optimizer^ Adafactor

Baselines:

- (1) T5 model tuning, default
- (2) T5 multitask

+ сравнение с GPT-3 few-shot

Результаты

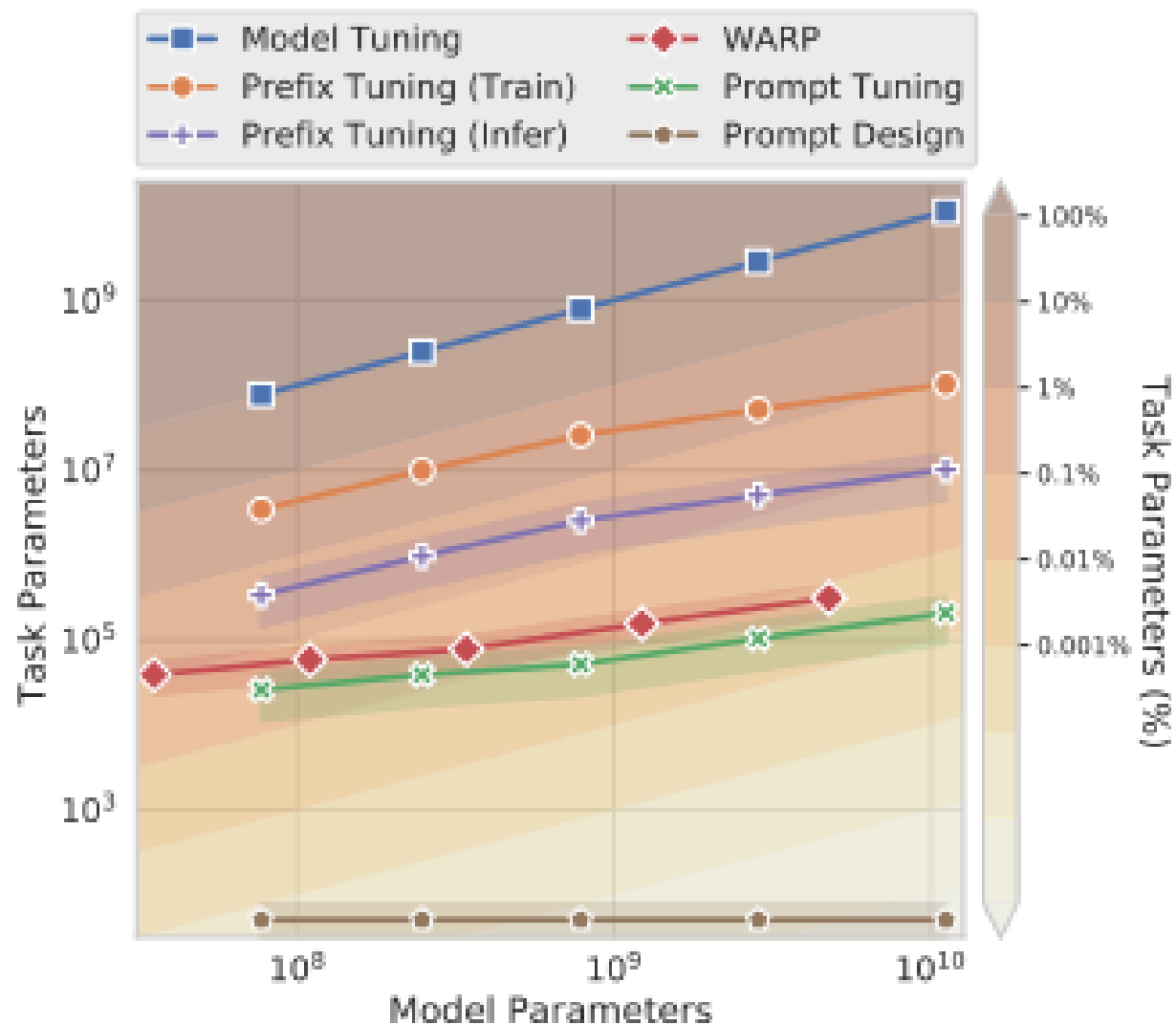
Model Tuning: все параметры специфичны для задания

Prefix Tuning: активации в prefix каждого слоя, 0.1–1% task-specific параметров для задачи, но больше для обучения

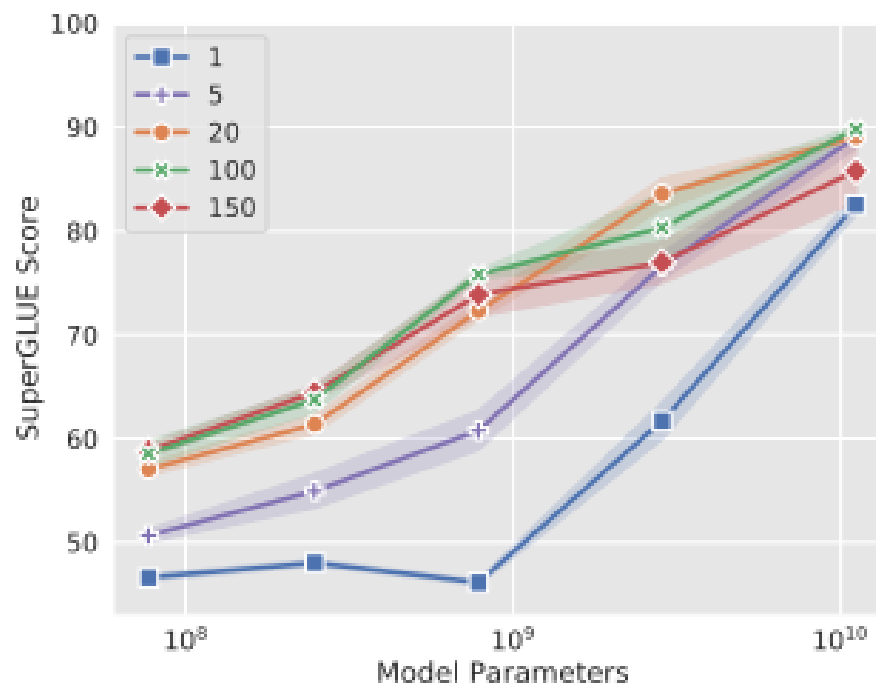
WARP: 0.1%, tuning входных/выходных слоев

Prompt Tuning: только prompt embeddings

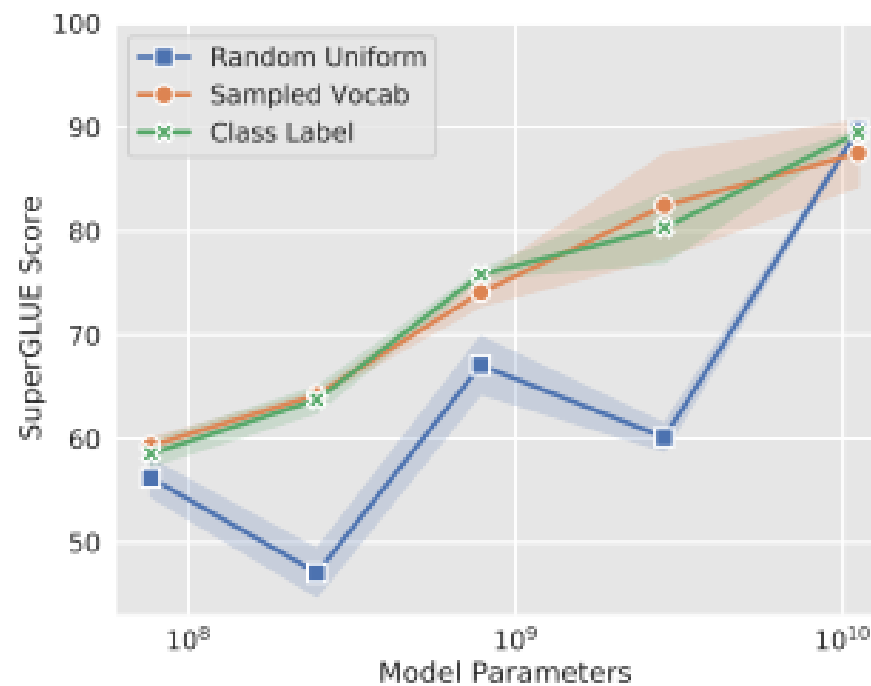
Prompt Design: только последовательность prompt IDs (500–2000 токенов)



Результаты

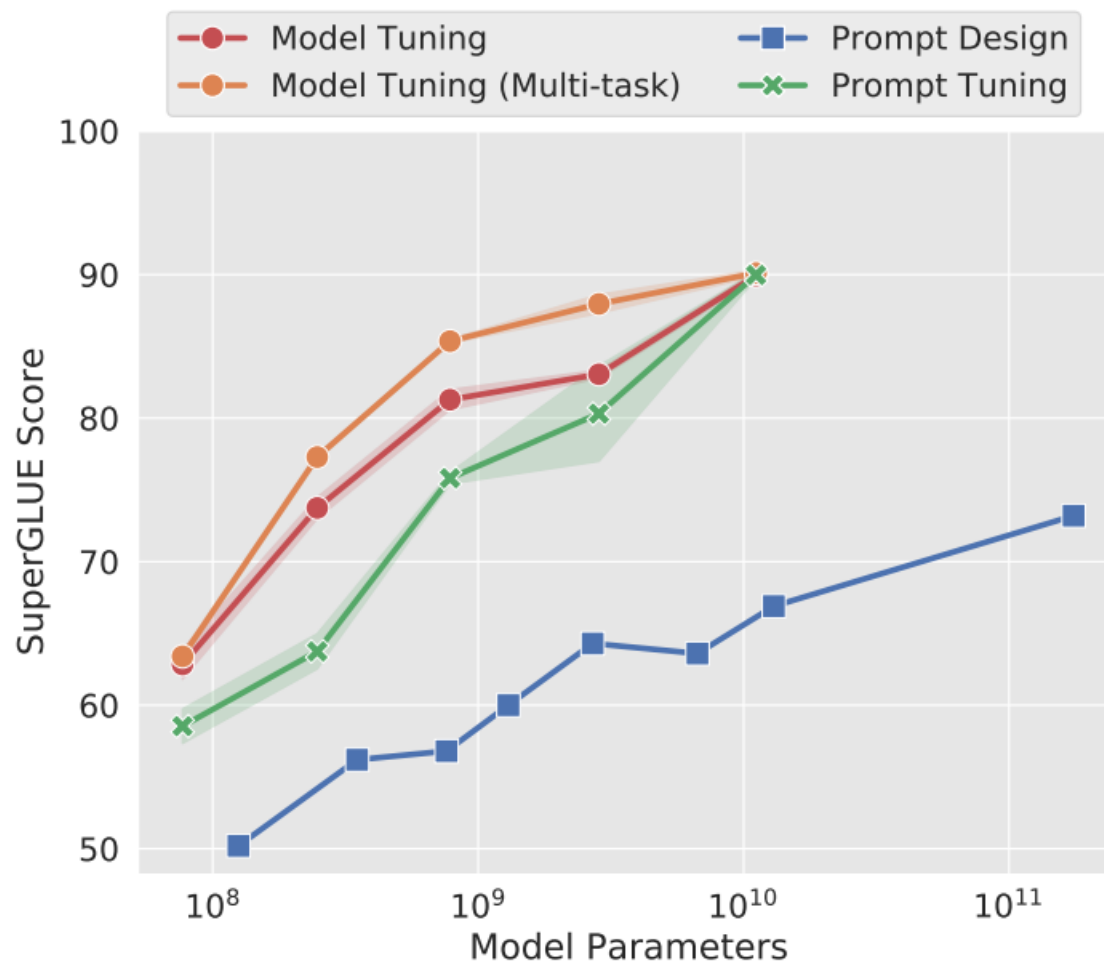


(a) Prompt length



(b) Prompt initialization

Результаты



Стандартная T5 - высокая производительность, но требует хранения отдельных копий модели для каждой конечной задачи.

Prompt tuned T5 повторно использует одну модель для всех задач.

Подход значительно превосходит несколько вариантов fewshot prompt для GPT-3.

Показано среднее значение и стандартное отклонение за 3 прогона для tuning'a

Результаты: ансамблирование

Dataset	Metric	Average	Best	Ensemble
BoolQ	acc.	91.1	91.3	91.7
CB	acc./F1	99.3 / 99.0	100.00 / 100.00	100.0 / 100.0
COPA	acc.	98.8	100.0	100.0
MultiRC	EM/F1 _a	65.7 / 88.7	66.3 / 89.0	67.1 / 89.4
ReCoRD	EM/F1	92.7 / 93.4	92.9 / 93.5	93.2 / 93.9
RTE	acc.	92.6	93.5	93.5
WiC	acc.	76.2	76.6	77.4
WSC	acc.	95.8	96.2	96.2
SuperGLUE (dev)		90.5	91.0	91.3

Table 3: Performance of a five-prompt ensemble built from a single frozen T5-XXL model exceeds both the average and the best among the five prompts.

PPT подход

PT ~ full model tuning

!

PT << few-shot learning

PPT подход

PT ~ full model tuning

!

PT << few-shot learning

Решение: предобучать soft prompts,
добавив их в стадию предобучения

PPT подход

PT ~ full model tuning

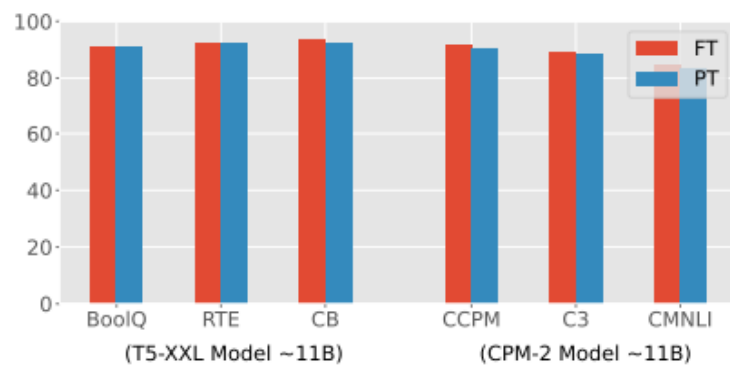
!

PT << few-shot learning

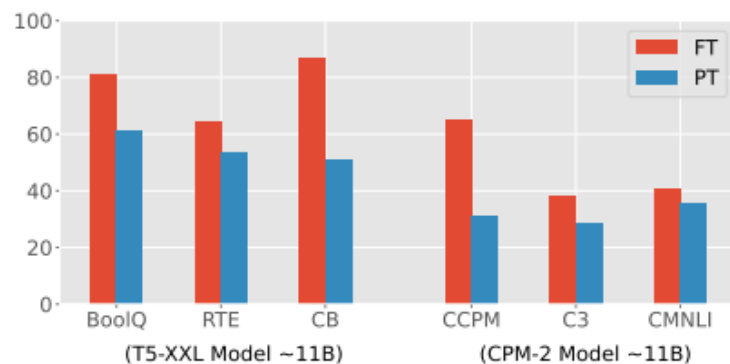
Решение: предобучать soft prompts,
добавив их в стадию предобучения

+ для обобщаемости похожие задачи классификации сводить в
единую и предобучать под неё

PPT подход



(a) Full-Data



(b) Few-Shot

Figure 2: Comparison between PT and FT. The tuned prompt is composed of 100 learnable embeddings whose dimensions are the same as the token embeddings of PLMs (4096 dimensions). All these results are based on 11B PLMs T5 and CPM-2. FT needs to optimize all 11B parameters, while PT only trains about 410K prompt parameters.

RPT подход

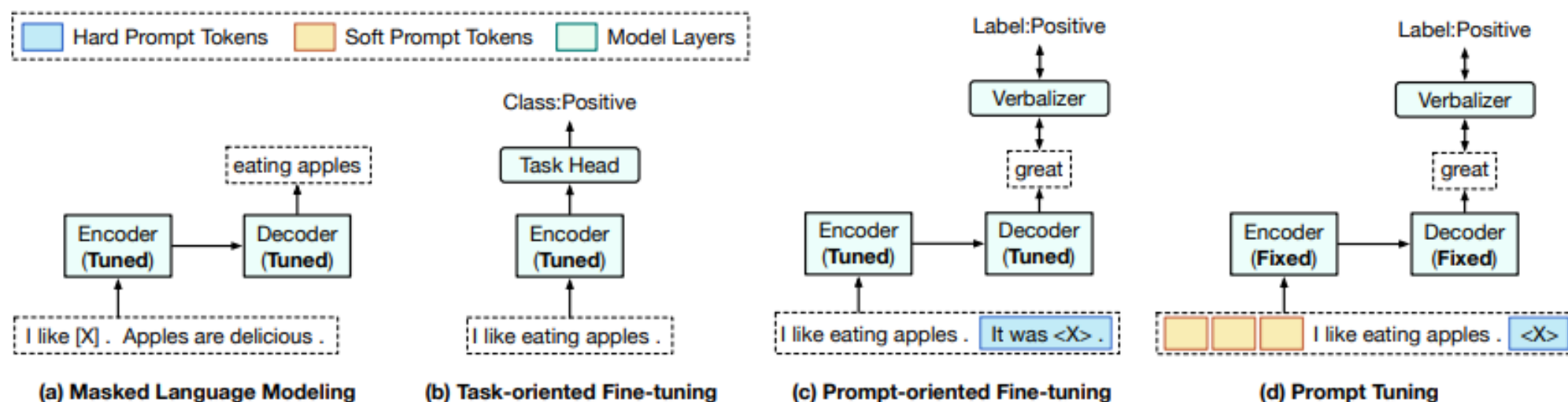


Figure 1: Paradigms of pre-training (masked language modeling), full-model tuning (task-oriented fine-tuning and prompt-oriented fine-tuning), and prompt tuning. The verbalizer is a function to map task labels to concrete words. $\langle X \rangle$ means the mask of typical pre-trained encoder-decoder models

PPT подход

$\{T_1, T_2, \dots, T_m\}$ – набор задач

$$T_i = \{PVP_i^1, PVP_i^2, \dots, PVP_i^{n_i}\}$$

PPT подход

$\{T_1, T_2, \dots, T_m\}$ – набор задач

$$T_i = \{PVP_i^1, PVP_i^2, \dots, PVP_i^{n_i}\}$$

$$PVP_i^k = (f_i^k, v_i^k)$$

f – mapping pattern, v - vocab

PPT подход

$\{T_1, T_2, \dots, T_m\}$ – набор задач

$$T_i = \{PVP_i^1, PVP_i^2, \dots, PVP_i^{n_i}\}$$

$$PVP_i^k = (f_i^k, v_i^k)$$

f – mapping pattern, v - vocab

$$PVP_i^{pre} = (f_i^{pre}, v_i^{pre})$$

PPT подход

$\{T_1, T_2, \dots, T_m\}$ – набор задач

$$T_i = \{PVP_i^1, PVP_i^2, \dots, PVP_i^{n_i}\}$$

$$PVP_i^k = (f_i^k, v_i^k)$$

f – mapping pattern, v - vocab

$$PVP_i^{pre} = (f_i^{pre}, v_i^{pre})$$

$\{P_1, P_2, \dots, P_m\}$ - prompts

PPT подход: setup для экспериментов

Chinese and English

PT on CPM-2 vs T5

PT on T5-XXL vs T5

PPT подход: setup для экспериментов

Chinese and English

PT on CPM-2 vs T5

PT on T5-XXL vs T5

- 100 soft tokens
- 410K params in PT (vs 11B FT)

РРТ
ПОДХОД:
ИТОГИ

English Tasks									
	Model	Method	SST-2 Acc.	SST-5 Acc.	RACE-m Acc.	RACE-h Acc.	BoolQ Acc.	RTE Acc.	CB F1
FT (11B)	T5-Small	-	72.8 _{3.1}	31.1 _{0.4}	26.4 _{0.6}	26.3 _{0.5}	59.2 _{0.6}	54.0 _{1.7}	70.1 _{4.6}
	T5-Base	-	74.6 _{2.7}	28.8 _{1.8}	27.2 _{0.5}	26.7 _{0.2}	61.9 _{2.1}	56.1 _{2.3}	70.4 _{2.6}
	T5-Large	-	89.1 _{2.2}	42.4 _{1.2}	48.2 _{1.6}	43.2 _{1.7}	74.6 _{0.9}	64.4 _{3.4}	82.3 _{2.2}
	T5-XL	-	89.6 _{3.2}	38.4 _{5.1}	55.0 _{2.8}	50.9 _{2.6}	77.2 _{2.1}	62.3 _{6.8}	81.9 _{9.0}
	T5-XXL	-	91.4 _{0.8}	40.6 _{2.0}	62.9_{3.9}	54.8_{3.0}	80.8 _{2.4}	64.1 _{2.0}	86.5_{5.3}
PT (410K)	T5-XXL	Vanilla PT	70.5 _{15.5}	32.3 _{8.3}	34.7 _{8.2}	31.6 _{3.5}	61.0 _{5.3}	53.5 _{3.5}	50.7 _{4.1}
		Hybrid PT	87.6 _{6.6}	40.9 _{2.7}	53.5 _{8.2}	44.2 _{6.4}	79.8 _{1.5}	56.8 _{2.6}	66.5 _{7.2}
		LM Adaption	77.6 _{7.5}	36.2 _{3.6}	27.3 _{0.2}	26.5 _{0.4}	62.0 _{0.3}	55.3 _{1.0}	61.2 _{1.7}
		PPT	93.5 _{0.3}	50.2_{0.7}	60.0 _{1.2}	53.0 _{0.4}	66.4 _{5.7}	58.9 _{1.6}	71.2 _{6.2}
		Hybrid PPT	93.8 _{0.1}	50.1 _{0.5}	62.5 _{0.9}	52.2 _{0.7}	82.0_{1.0}	59.8 _{3.2}	73.2 _{7.0}
		Unified PPT	94.4_{0.3}	46.0 _{1.3}	58.0 _{0.9}	49.9 _{1.3}	76.0 _{2.7}	65.8_{2.1}	82.2 _{5.4}
Chinese Tasks									
	Model	Method	ChnSent Acc.	Amazon Acc.	CCPM Acc.	C ³ Acc.	LCQMC Acc.	CMNLI Acc.	OCNLI Acc.
FT (11B)	mT5-Small	-	76.1 _{2.6}	29.9 _{1.9}	31.9 _{1.2}	29.6 _{0.5}	52.4 _{2.5}	36.5 _{0.2}	34.9 _{1.3}
	mT5-Base	-	78.2 _{0.6}	36.4 _{0.9}	40.4 _{6.8}	29.4 _{0.6}	50.9 _{1.0}	36.3 _{0.5}	35.4 _{0.6}
	mT5-Large	-	79.1 _{0.6}	31.0 _{1.4}	46.0 _{4.0}	29.9 _{0.8}	52.1 _{0.6}	35.8 _{1.2}	35.2 _{1.1}
	mT5-XL	-	82.7 _{2.6}	35.5 _{1.7}	68.3 _{5.1}	29.7 _{1.2}	52.9 _{2.4}	36.8 _{1.6}	35.6 _{0.5}
	mT5-XXL	-	83.6 _{1.5}	42.1 _{0.8}	79.7 _{1.1}	37.2 _{3.3}	53.1 _{1.0}	39.0 _{0.4}	37.4 _{1.2}
	CPM-2	-	86.1 _{1.8}	42.5 _{2.0}	81.8 _{1.6}	38.4 _{3.7}	58.8 _{1.8}	40.7 _{1.0}	38.5 _{1.5}
PT (410K)	CPM-2	Vanilla PT	62.1 _{3.1}	30.3 _{4.8}	31.0 _{9.7}	28.2 _{0.4}	51.5 _{3.4}	35.4 _{0.5}	37.0 _{0.5}
		Hybrid PT	79.2 _{4.0}	39.1 _{3.8}	46.6 _{15.0}	29.2 _{0.5}	54.6 _{2.3}	37.1 _{0.6}	37.8 _{1.4}
		LM Adaption	74.3 _{5.2}	35.2 _{2.4}	33.7 _{12.8}	30.2 _{1.5}	51.4 _{2.9}	35.1 _{0.3}	38.0 _{1.1}
		PPT	90.1 _{0.8}	48.6 _{0.6}	85.4_{0.6}	43.8 _{2.2}	59.1 _{0.6}	43.0_{0.5}	40.1 _{0.4}
		Hybrid PPT	89.5 _{0.3}	48.8_{2.0}	83.9 _{0.5}	46.0 _{0.5}	67.3_{0.9}	41.3 _{0.8}	38.7 _{0.6}
		Unified PPT	90.7_{0.2}	44.6 _{1.1}	83.4 _{0.9}	50.2_{0.6}	55.0 _{0.4}	40.6 _{0.4}	41.5_{1.5}