

# Recycling GPT-2 and Don't stop pretraining

Feb 17 2021

# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- Problem
  - Training GPT-2 from scratch is data and computation intensive
  - Can we adapt English GPT-2 in an efficient way to genetically related languages such as Italian and Dutch(same word order SVO, adj-noun)?
- Solution
  - Retrain lexical embeddings of English GPT-2 on Italian/Dutch data without tuning the transformer layers
  - Bonus – Lexical embeddings for Italian/Dutch are aligned with English and can induce a bilingual lexicon
- Caveat – Training a medium/large size model is still time-consuming
  - Transform relearned lexical embeddings of GPT-2 small to GPT-2 medium
- Results – GPT-2 with relearned lexical embeddings on Italian can generate realistic sentences. Gets better with some finetuning.

# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- GPT-2
  - Auto-regressive Transformer-decoder base language model for English
  - 4 sizes: **small (sml)**, **medium (med)**, large, extra large
- Pre-training data
  - Italian
    - Same pre-training data as Italian GPT-2 (GepPpeTto) – Wiki (2.8GB) + ItWac (11GB)
  - Dutch
    - Wiki (2GB) + news (2.9GB) + books (6.5GB) + Dutch news (2.1GB)
- Evaluation corpora (e.g., perplexity)
  - Italian - Eval corpora from GepPpeTto
  - Dutch – 22-genre, 500M work SoNaR

# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- Tokenization
  - Dutch – full pre-training data of GePpeTto
    - Vocab size – balance of large vocabulary vs coverage for uncommon tokens
- Computation
  - 8 x V100 32GB GPU
  - 16-bit automatic mixed-precision training (2-3 times reduction)
  - Split each document into a windows of 128 tokens (for retraining lexical embeddings)
  - Bucketed random sampling (decreases amount of padding within minibatches)
  - Maximum batch size that fit into GPU memory

# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- Cross-language transfer
  - English GPT-2 uses English BPE vocabulary, not usable for new languages
  - Randomly initialize the lexical embedding layer (first and last layer) + freeze rest of the model + relearn the lexical embeddings on language specific data
- Relearning lexical embeddings
  - Still sample training time but converges faster
  - Sml > med in terms of perplexity
- Vocabulary alignment
  - Relearned Italian/Dutch lexical embeddings are aligned to English embeddings

Model	PPL	
	ita	nld
sml <sub>rle</sub>	44.19	48.85
med <sub>rle</sub> *	-	185.02

Table 1: Perplexity scores on test data for relearned lexical embeddings (rle). \* med training from random initialisation is stopped early after two days of training.

English	Italian	Dutch
while	mentre	terwijl
genes	geni	genen
clothes	vestiti	kleren
musicians	composi[...]	artiesten
permitted	ammessa	toegelaten
Finally	infine	Eindelijk
satisfied	soddisfatto	tevreden
Accuracy:	85%	89%

Table 2: Alignment of closest tokens in the lexical embeddings of sml<sub>rle</sub> for Italian and Dutch. Accuracy scores are based on a manual evaluation by the authors of 200 random aligned tokens. Semantically correct subword matches are included.

# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- Text generation
  - Sometime English word order, doesn't produce correct Italian prepositional articles. Get's noun-adjective order for Italian
  - Syntactic features of tokens in addition to semantics

Italian	Literal English translation
La prima parte del film venne <i>distribuito</i> in Giappone con l'aggiunta della colonna sonora.	The first part of the film was <i>distributed</i> in Japan with the addition of the soundtrack.
L'unico motivo <i>di la</i> mia insoddisfazione fu il fatto che l'inizio della sua attività [...]	The only reason <i>of the</i> my unsatisfaction was the fact that the beginning of-the his/her activity [...]
Il suo nome deriva da un vocabolo arabo.	The his/her name derives from a word Arabic.
Dutch	Literal English translation
In een artikel in de Journal of Economicologie (1998), <i>The New York Times schrijft</i> :	In an article in the Journal of Economicology (1998), <i>The New York Times writes</i> :
Ik kan me niet voorstellen dat mensen van mijn generatie <i>zijn zo boos op mij te wachten</i> .	I can me not imagine that people of my generation <i>are so mad at me to wait</i> .
Ik heb niets gedaan om mijn moeder te helpen.	I have nothing done to my mother to help.

Table 3: A selection of generated sentences by the `sml` model with Italian and Dutch lexical embeddings. Parts in Italic are ungrammatical in the target language.

# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- Scaling up complexity
  - Relearning for large models are still time consuming
  - Idea: Use the vocabulary alignment between source and target languages for smaller model can be used to initialize embeddings of larger model
- Approach
  - 50K English words
  - Train: Learn a mapping between a English word from small model to medium model
  - Test: Map the Italian words from small model to medium model
- Learning mapping
  - Regression – Learn  $W$  that minimizes Euclidean distance between source and target embedding

# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- Orthogonal Procrustes

- Constrain the transformation  $W$  to be orthogonal matrix
- Exact solution, only rotations, preserves monolingual invariance

- Weighted K-NN

- Unknown target language token = k-nn source embedding space for English and use the distance-weighted sum of these tokens in the target space

Model	Italian		Dutch		
	Int@1k	PPL	Int@1k	PPL	PPL (1 epoch)
med <sub>r1e</sub> (1 epoch)	0.38	-	185.02	-	-
sml <sub>r1e</sub> $\xrightarrow{proc}$ med	<b>0.61</b>	$8.12 \times 10^{12}$	<b>0.61</b>	$5.02 \times 10^{12}$	52.69
sml <sub>r1e</sub> $\xrightarrow{lstsq}$ med	0.56	<b>364.06</b>	0.56	<b>293.61</b>	<b>47.57</b>
sml <sub>r1e</sub> $\xrightarrow{1-nn}$ med	0.37	2,764.19	0.36	1,101.59	50.25
sml <sub>r1e</sub> $\xrightarrow{10-nn}$ med	0.37	20,715.80	0.35	11,871.66	56.88

Table 4: Scores for different transformation methods. Int@1K are the average 1k nearest English neighbours intersection (int) fractions between sml and transformed med embeddings. PPL is the perplexity on the test sets for Italian and Dutch. PPL (1 epoch) indicates the perplexity after one epoch of training, indicating closeness of the transformation to a good local optimum.



# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- Full model finetuning
  - Sml model - relearn + finetune
  - Med model – relearn + align + finetune
- Automatic Evaluation
  - sml + finetuning is best

Model	PPL	
	ita	nld
sml <sub>rle</sub>	44.19	48.85
sml <sub>rle</sub> + finetuning	<b>42.45</b>	<b>39.59</b>
med <sub>rle</sub>	42.51	44.68
GePpeTto <sub>sml</sub>	106.84	-

Table 5: Perplexities of the concatenated test data for the final models. The med<sub>rle</sub> model is in practice the sml<sub>rle</sub>  $\xrightarrow{lstsq}$  med model. The small Dutch model seems to benefit more from full model finetuning.

Model	Proceedings	News	Legal
sml <sub>rle</sub>	44.47	239.14	52.01
sml <sub>rle</sub> + fine	<b>36.35</b>	<b>171.83</b>	<b>42.92</b>
med <sub>rle</sub>	40.62	234.52	45.01

Table 7: Perplexities for some SoNaR genres in Dutch. Models rankings are consistent across genres.

Model	Social	News	Legal
sml <sub>rle</sub>	134.64	67.14	16.95
sml <sub>rle</sub> + fine	<b>118.19</b>	<b>55.63</b>	15.36
med <sub>rle</sub>	123.64	59.18	<b>14.95</b>
GePpeTto <sub>sml</sub>	179.47	80.83	34.71

Table 6: Perplexities for different genres within the Italian test data. Rankings are consistent with Table 5 except for the legal domain.

# As good as new. How to successfully recycle English GPT-2 to make models for other languages (Vries & Nissim, 2020)

- Human Evaluation
  - Generations from relearn are easily identifiable
  - Generations from relearn + finetune are not so.

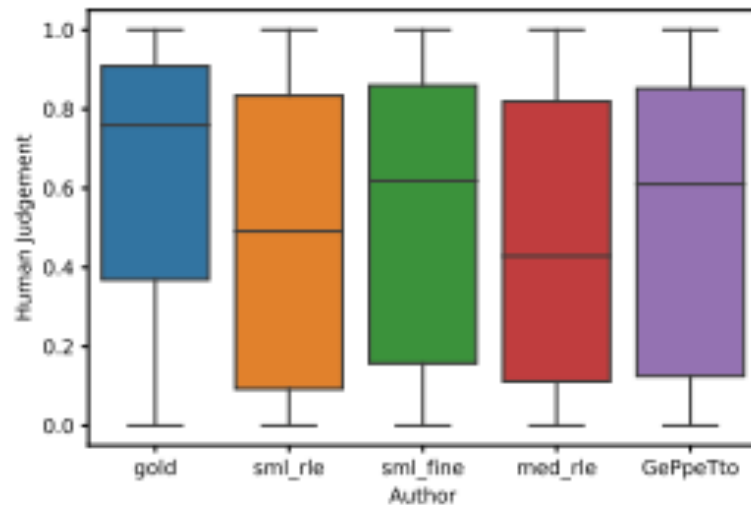


Figure 1: Human judgement scores for Italian texts.

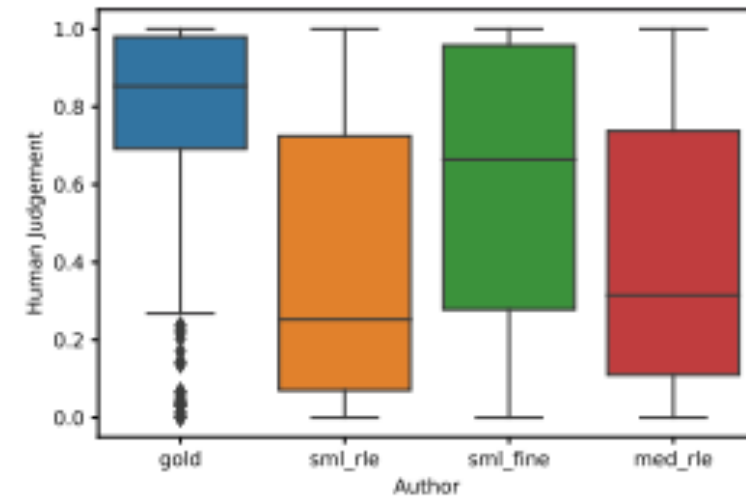


Figure 2: Human judgement scores for Dutch texts.

# Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., ACL'20)

- Problem
  - Do latest pretrained LM work universally or is it still helpful to build separate pretrained models for specific domains?
  - How benefits of continued pretraining varies with factors like amount of labeled task data, proximity of target domain to original pretrain corpus?
- Results
  - RoBERTa + pretrain (domain unlabeled) improves for both low/high resource settings (Domain adaptive pretraining - DAPT)
  - RoBERTa + pretrain (task unlabeled) improves with or without DAPT (Task adaptive pretraining)
  - RoBERTa + pretrain (automatically selected task specific unlabeled) cost effective approach

# Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., ACL'20)

- Domains
  - Biomedical (BIOMED)
  - Computer science (CS)
  - Newstext (RealNews)
  - Amazon reviews (Amazon)
- Vocab overlap (top 10K most freq. unigrams)
- Expectation: More dissimilar the domain from the pretraining domain, higher the potential for DAPT

PT	100.0	54.1	34.5	27.3	19.2
News	54.1	100.0	40.0	24.9	17.3
Reviews	34.5	40.0	100.0	18.3	12.7
BioMed	27.3	24.9	18.3	100.0	21.4
CS	19.2	17.3	12.7	21.4	100.0
	PT	News	Reviews	BioMed	CS

Figure 2: Vocabulary overlap (%) between domains. PT denotes a sample from sources similar to ROBERTA's pretraining corpus. Vocabularies for each domain are created by considering the top 10K most frequent words (excluding stopwords) in documents sampled from each domain.

# Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., ACL'20)

- DAPT MLM loss

Domain	Pretraining Corpus	# Tokens	Size	$\mathcal{L}_{\text{ROB.}}$	$\mathcal{L}_{\text{DAPT}}$
BIOMED	2.68M full-text papers from S2ORC (Lo et al., 2020)	7.55B	47GB	1.32	0.99
CS	2.22M full-text papers from S2ORC (Lo et al., 2020)	8.10B	48GB	1.63	1.34
NEWS	11.90M articles from REALNEWS (Zellers et al., 2019)	6.66B	39GB	1.08	1.16
REVIEWS	24.75M AMAZON reviews (He and McAuley, 2016)	2.11B	11GB	2.10	1.93
ROBERTA (baseline)	see Appendix §A.1	N/A	160GB	<sup>‡</sup> 1.19	-

Table 1: List of the domain-specific unlabeled datasets. In columns 5 and 6, we report ROBERTA’s masked LM loss on 50K randomly sampled held-out documents from each domain before ( $\mathcal{L}_{\text{ROB.}}$ ) and after ( $\mathcal{L}_{\text{DAPT}}$ ) DAPT (lower implies a better fit on the sample). <sup>‡</sup> indicates that the masked LM loss is estimated on data sampled from sources *similar* to ROBERTA’s pretraining corpus.

# Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., ACL'20)

- Downstream classification tasks

Domain	Task	Label Type	Train (Lab.)	Train (Unl.)	Dev.	Test	Classes
BIOMED	CHEMPROT	relation classification	4169	-	2427	3469	13
	<sup>†</sup> RCT	abstract sent. roles	18040	-	30212	30135	5
CS	ACL-ARC	citation intent	1688	-	114	139	6
	SciERC	relation classification	3219	-	455	974	7
NEWS	HYPERPARTISAN	partisanship	515	5000	65	65	2
	<sup>†</sup> AGNEWS	topic	115000	-	5000	7600	4
REVIEWS	<sup>†</sup> HELPFULNESS	review helpfulness	115251	-	5000	25000	2
	<sup>†</sup> IMDB	review sentiment	20000	50000	5000	25000	2

Table 2: Specifications of the various target task datasets. <sup>†</sup> indicates high-resource settings. Sources: CHEMPROT (Kringelum et al., 2016), RCT (Dernoncourt and Lee, 2017), ACL-ARC (Jurgens et al., 2018), SciERC (Luan et al., 2018), HYPERPARTISAN (Kiesel et al., 2019), AGNEWS (Zhang et al., 2015), HELPFULNESS (McAuley et al., 2015), IMDB (Maas et al., 2011).

Dom.	Task	ROBA.	DAPT	$\neg$ DAPT
BM	CHEMPROT	81.9 <sub>1.0</sub>	<b>84.2</b> <sub>0.2</sub>	79.4 <sub>1.3</sub>
	<sup>†</sup> RCT	87.2 <sub>0.1</sub>	<b>87.6</b> <sub>0.1</sub>	86.9 <sub>0.1</sub>
CS	ACL-ARC	63.0 <sub>5.8</sub>	<b>75.4</b> <sub>2.5</sub>	66.4 <sub>4.1</sub>
	SciERC	77.3 <sub>1.9</sub>	<b>80.8</b> <sub>1.5</sub>	79.2 <sub>0.9</sub>
NEWS	HYP.	86.6 <sub>0.9</sub>	<b>88.2</b> <sub>5.9</sub>	76.4 <sub>4.9</sub>
	<sup>†</sup> AGNEWS	<b>93.9</b> <sub>0.2</sub>	<b>93.9</b> <sub>0.2</sub>	93.5 <sub>0.2</sub>
REV.	<sup>†</sup> HELPFUL.	65.1 <sub>3.4</sub>	<b>66.5</b> <sub>1.4</sub>	65.1 <sub>2.8</sub>
	<sup>†</sup> IMDB	95.0 <sub>0.2</sub>	<b>95.4</b> <sub>0.2</sub>	94.1 <sub>0.4</sub>

Table 3: Comparison of ROBERTA (ROBA.) and DAPT to adaptation to an *irrelevant* domain ( $\neg$ DAPT). Reported results are test macro- $F_1$ , except for CHEMPROT and RCT, for which we report micro- $F_1$ , following Beltagy et al. (2019). We report averages across five random seeds, with standard deviations as subscripts. <sup>†</sup> indicates high-resource settings. Best task performance is boldfaced. See §3.3 for our choice of irrelevant domains.



# Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., ACL'20)

- Task adaptive pretraining
  - Datasets curated to capture specific tasks of interest tend to cover only a subset of text available within the broader domain
  - CHEMPROT – extracting relations between chemicals and proteins focusing on abstracts from hand selected PubMed articles
  - Hypothesis – Cases where task data is narrowly defined subset of broader domain, pretraining on task dataset itself or data relevant to the task may be helpful
  - TAPT is faster than DAPT

Domain	Task	RoBERTA	Additional Pretraining Phases		
			DAPT	TAPT	DAPT + TAPT
BioMed	CHEMPROT	81.9 <sub>1.0</sub>	84.2 <sub>0.2</sub>	82.6 <sub>0.4</sub>	<b>84.4</b> <sub>0.4</sub>
	†RCT	87.2 <sub>0.1</sub>	87.6 <sub>0.1</sub>	87.7 <sub>0.1</sub>	<b>87.8</b> <sub>0.1</sub>
CS	ACL-ARC	63.0 <sub>5.8</sub>	75.4 <sub>2.5</sub>	67.4 <sub>1.8</sub>	<b>75.6</b> <sub>3.8</sub>
	SciERC	77.3 <sub>1.9</sub>	80.8 <sub>1.5</sub>	79.3 <sub>1.5</sub>	<b>81.3</b> <sub>1.8</sub>
NEWS	HYPERPARTISAN	86.6 <sub>0.9</sub>	88.2 <sub>5.9</sub>	<b>90.4</b> <sub>5.2</sub>	90.0 <sub>6.6</sub>
	†AGNEWS	93.9 <sub>0.2</sub>	93.9 <sub>0.2</sub>	94.5 <sub>0.1</sub>	<b>94.6</b> <sub>0.1</sub>
REVIEWS	†HELPPFULNESS	65.1 <sub>3.4</sub>	66.5 <sub>1.4</sub>	68.5 <sub>1.9</sub>	<b>68.7</b> <sub>1.8</sub>
	†IMDB	95.0 <sub>0.2</sub>	95.4 <sub>0.1</sub>	95.5 <sub>0.1</sub>	<b>95.6</b> <sub>0.1</sub>

Table 5: Results on different phases of adaptive pretraining compared to the baseline RoBERTA (col. 1). Our approaches are DAPT (col. 2, §3), TAPT (col. 3, §4), and a combination of both (col. 4). Reported results follow the same format as Table 3. State-of-the-art results we can compare to: CHEMPROT (84.6), RCT (92.9), ACL-ARC (71.0), SciERC (81.8), HYPERPARTISAN (94.8), AGNEWS (95.5), IMDB (96.2); references in §A.2.

# Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., ACL'20)

- Task adaptive pretraining
  - TAPT optimizes for single task performance, to the detriment of cross-task transfer

BIOMED	RCT	CHEMPROT	CS	ACL-ARC	SCIERC
TAPT	87.7 <sub>0.1</sub>	82.6 <sub>0.5</sub>	TAPT	67.4 <sub>1.8</sub>	79.3 <sub>1.5</sub>
Transfer-TAPT	87.1 <sub>0.4</sub> (↓0.6)	80.4 <sub>0.6</sub> (↓2.2)	Transfer-TAPT	64.1 <sub>2.7</sub> (↓3.3)	79.1 <sub>2.5</sub> (↓0.2)
NEWS	HYPERPARTISAN	AGNEWS	REVIEWS	HELPFULNESS	IMDB
TAPT	89.9 <sub>9.5</sub>	94.5 <sub>0.1</sub>	TAPT	68.5 <sub>1.9</sub>	95.7 <sub>0.1</sub>
Transfer-TAPT	82.2 <sub>7.7</sub> (↓7.7)	93.9 <sub>0.2</sub> (↓0.6)	Transfer-TAPT	65.0 <sub>2.6</sub> (↓3.5)	95.0 <sub>0.1</sub> (↓0.7)

Table 6: Though TAPT is effective (Table 5), it is harmful when applied *across* tasks. These findings illustrate differences in task distributions within a domain.



# Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., ACL'20)

- Curated TAPT
  - Large unlabeled corpus + downsampled to collect task annotations
  - Large unlabeled corpus should have similar distribution to task's training data
  - Data
    - Downsample RCT to 500 (180K unlabeled)
  - DAPT + Curated TAPT is best
  - Curating large amounts of data from task distribution is extremely beneficial

Pretraining	BIOMED RCT-500	NEWS HYP.	REVIEWS IMDB <sup>†</sup>
TAPT	79.8 <sub>1.4</sub>	90.4 <sub>5.2</sub>	95.5 <sub>0.1</sub>
DAPT + TAPT	83.0 <sub>0.3</sub>	90.0 <sub>6.6</sub>	95.6 <sub>0.1</sub>
Curated-TAPT	83.4 <sub>0.3</sub>	89.9 <sub>9.5</sub>	95.7 <sub>0.1</sub>
DAPT + Curated-TAPT	<b>83.8</b> <sub>0.5</sub>	<b>92.1</b> <sub>3.6</sub>	<b>95.8</b> <sub>0.1</sub>

Table 7: Mean test set macro- $F_1$  (for HYP. and IMDB) and micro- $F_1$  (for RCT-500), with Curated-TAPT across five random seeds, with standard deviations as subscripts. <sup>†</sup> indicates high-resource settings.

# Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., ACL'20)

- Automated Data Selection for TAPT
  - Embed task-relevant data from domain and from task in shared space
  - Select candidates from the domain based on queries using the task data
  - VAMPIRE, BoW LM
  - kNN-TAPT
  - RAND-TAPT

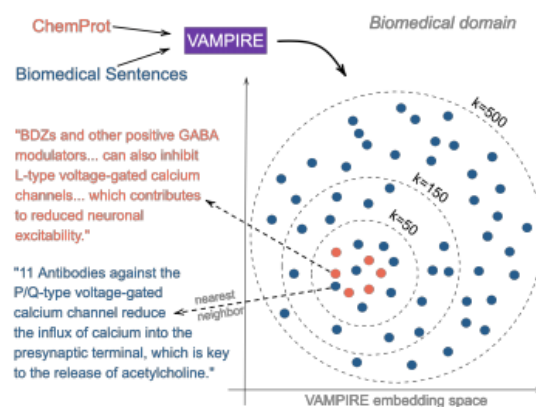


Figure 3: An illustration of automated data selection (§5.2). We map unlabeled CHEMPROT and 1M BIOMED sentences to a shared vector space using the VAMPIRE model trained on these sentences. Then, for each CHEMPROT sentence, we identify  $k$  nearest neighbors, from the BIOMED domain.

Pretraining	BIOMED		CS
	CHEMPROT	RCT-500	ACL-ARC
ROBERTA	81.9 <sub>1.0</sub>	79.3 <sub>0.6</sub>	63.0 <sub>5.8</sub>
TAPT	82.6 <sub>0.4</sub>	79.8 <sub>1.4</sub>	67.4 <sub>1.8</sub>
RAND-TAPT	81.9 <sub>0.6</sub>	80.6 <sub>0.4</sub>	69.7 <sub>3.4</sub>
50NN-TAPT	83.3 <sub>0.7</sub>	80.8 <sub>0.6</sub>	70.7 <sub>2.8</sub>
150NN-TAPT	83.2 <sub>0.6</sub>	81.2 <sub>0.8</sub>	73.3 <sub>2.7</sub>
500NN-TAPT	83.3 <sub>0.7</sub>	81.7 <sub>0.4</sub>	<b>75.5<sub>1.9</sub></b>
DAPT	<b>84.2<sub>0.2</sub></b>	<b>82.5<sub>0.5</sub></b>	75.4 <sub>2.5</sub>

Table 8: Mean test set micro- $F_1$  (for CHEMPROT and RCT) and macro- $F_1$  (for ACL-ARC), across five random seeds, with standard deviations as subscripts, comparing RAND-TAPT (with 50 candidates) and  $k$ NN-TAPT selection. Neighbors of the task data are selected from the domain data.