

RADS: Reinforcement Learning-Based Sample Selection Improves Transfer Learning in Low-resource and Imbalanced Clinical Settings

Anonymous ACL submission

Abstract

A common strategy in transfer learning is few shot fine-tuning, but its success is highly dependent on the quality of samples selected as training examples. Active learning methods such as uncertainty sampling and diversity sampling can select useful samples. However, under extremely low-resource and class-imbalanced conditions, they often favor outliers rather than truly informative samples, resulting in degraded performance. In this paper, we introduce **RADS (Reinforcement Domain Adaptive Sampling)**, a robust sample selection strategy using reinforcement learning (RL) to identify the most informative samples. Experimental evaluations on several real world clinical datasets show our sample selection strategy enhances model transferability while maintaining robust performance under extreme class imbalance compared to traditional methods.

1 Introduction

Maximizing the utility of limited data is a crucial focus of Natural Language Processing (NLP) research in domains such as clinical texts where acquiring large amounts of gold standard data may be difficult due to data restrictions and the relative rarity of many disease conditions. The high cost of annotation in such highly specialized domains further limits availability of labeled data. Yet, the effectiveness of NLP techniques in healthcare heavily relies on the quality of annotated datasets, particularly because clinical data contains specialized symbols, abbreviations, and medical jargon (Touvron et al., 2023; Liu et al., 2024a).

Transfer Learning (TL) (Tan et al., 2018), in which knowledge learned from a task is reused to boost performance on a different but related (target) task, has shown effectiveness across various machine learning applications (Weiss et al., 2016) and opens new avenues for addressing low-resource scenarios. Previous works have attempted to leverage pretrained embeddings (Maimaiti et al., 2021)

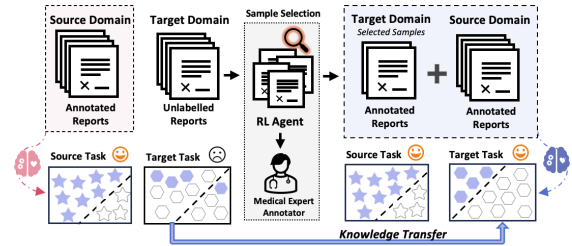


Figure 1: RL-based active sampling for transfer learning from source domain to target domain. Domain shift reduces zero-shot generalization from the source-trained model to the target domain. Our sample selection strategy uses RL to identify key samples from the target domain. By jointly fine-tuning on the selected target samples with the source data, the model achieves good performance on both domains.

and few shot examples (Alyafeai et al., 2020) to facilitate transfer learning in NLP. However, when the target task offers very few labeled instances, these approaches may generate unreliable outputs. This is an especially acute problem in healthcare, where reliability is paramount.

Class imbalance (Johnson and Khoshgoftaar, 2019) is another challenge for low-resource settings. In clinical datasets, there is often a scarcity of positive cases due to the low prevalence of many conditions, making such instances both highly valuable and limited in number. At the same time, differences in data collection protocols can lead some disease datasets to contain a very high proportion of positive samples. These extreme disparities in class distribution further hinder the transferability of NLP models across clinical datasets.

Clinical documentation is heterogeneous, reflecting diverse investigations including CT and PET scans, or cytology and histopathology analysis. Although disease detection cues appear across these document types, their content structures, terminology, and linguistic expressions can vary greatly. CT and PET scan reports primarily emphasize imaging-based findings (Townsend et al., 2004), whereas

068 cytology and histopathology reports focus on cellu- 118
069 lar and tissue-level observations (Jensen, 2021). 119

070 Previous works have shown the efficacy of NLP 120
071 techniques for disease detection from clinical re- 121
072 ports. However, models fine-tuned on one report 122
073 type show clear performance degradation when ap-
074 plied to another (Han et al., 2025). While disease-
075 related signals overlap to some extent across differ-
076 ent document types, existing disease detection mod-
077 els still fall short of human performance in transfer-
078 ring knowledge between them. As the preparation
079 of gold standard annotated datasets for training is
080 time-consuming, it is therefore important to ex-
081 plore effective knowledge transfer strategies from
082 existing datasets to new but similar tasks. This
083 not only improves annotation efficiency but also
084 enhances the models’ adaptability in dealing with
085 variations in task settings.

086 In this work, we propose RADS (Reinforcement 123
087 Domain Adaptive Sampling), a robust strategy for 124
088 knowledge transfer between related but distinct 125
089 sources. Following the active learning paradigm 126
090 (Fu et al., 2013), we enhance transfer learning 127
091 by identifying and selecting the most relevant sam- 128
092 ples for few-shot fine-tuning as shown in Figure 1. 129
093 First, we employ an RL-based agent to identify the 130
094 most informative samples within the target dataset. 131
095 These selected samples are then annotated by med- 132
096 ical experts and incorporated into the fine-tuning 133
097 process. By jointly fine-tuning the model on the 134
098 source dataset and the newly annotated target sam- 135
099 ples, the model is able to preserve strong perfor- 136
100 mance on the source domain while achieving im- 137
101 proved generalization to the target domain. We 138
102 evaluated this approach across multiple real world 139
103 clinical datasets. Experimental results show that 140
104 our method improves both the adaptability and per- 141
105 formance of disease detection between different 142
106 sources. In the context of transfer learning, this 143
107 technique offers a promising way to both reduce 144
108 annotation effort and enhance model robustness in 145
109 low-resource and class-imbalanced settings. 146

110 Our contributions are summarized as follows:

- 111 • This work is the first to address the challenges 147
112 posed by low-resource and class-imbalance 148
113 scenarios in disease detection with different 149
114 report types from real clinical data sources. 150
- 115 • We propose RADS, a robust RL-based sample 151
116 selection strategy tailored to scenarios with 152
117 both data scarcity and class imbalance. 153

- Extensive experiments on several clinical 118
datasets confirm that our transfer learning ap- 119
proach is more effective between similar but 120
different sources, even under low-resource 121
and class-imbalanced conditions. 122

2 Related Work 123

124 With high-quality annotated datasets, NLP methods 125
126 have shown promising results in disease detection. 127
128 Based on the concept features relevant to diseases, 129
130 dictionary-based detection approaches and classi- 131
132 cal machine learning have shown effective perfor- 133
134 mance (Rozova et al., 2023b; Martinez et al., 2015). 135
136 Bag-of-words models have also been utilized, of- 137
138 ten combined with machine learning techniques to 139
140 further enhance accuracy and scalability in disease 141
142 detection (Cury et al., 2021; López-Úbeda et al., 143
144 2020). Recently, large language models (LLMs), 145
146 such as BioBERT (Lee et al., 2020) and Clinical- 147
148 BERT (Huang et al., 2019), pre-trained on large 149
150 biomedical corpora, have improved contextual un- 151
152 derstanding in clinical texts (Consoli et al., 2024; 153
154 Han et al., 2025). 155

156 Low resource settings remain challenging for 157
158 NLP tasks. Few-shot fine-tuning (Brown et al., 159
160 2020; Gu et al., 2022; Liu et al., 2022), where large 161
162 pre-trained models are adapted using only a small 163
164 number of labeled examples, has shown promis- 164
165 ing results. Selecting effective few-shot samples 165
166 is critical, and active learning strategies such as 166
167 uncertainty sampling (Nguyen et al., 2022) and 167
168 diversity sampling (Yang et al., 2015) are often 168
169 employed. However, these methods typically optimize 169
170 a single metric, and under domain shift, tend to se- 170
171 lect distributional outliers rather than truly informa- 171
172 tive samples (Gonsior et al., 2024). Reinforcement 172
173 Learning (RL) (Fang et al., 2017; Liu et al., 2024b) 173
174 offers a potential solution by optimizing more flex- 174
175 ible and adaptive sample selection policies, thereby 175
176 improving robustness in different contexts. 176

177 Class imbalance is especially crucial in low- 177
178 resource clinical NLP tasks (Ghosh et al., 2024). 178
179 Data-level approaches, such as oversampling mi- 179
180 nority classes (Hairani et al., 2024) and undersam- 180
181 pling majority classes (Yang et al., 2024), are typi- 181
182 cally used to balance class distributions. Algorithm- 182
183 level methods, such as cost-sensitive learning (Araf 183
184 et al., 2024) and focal loss adjustments (Aljohani 184
185 et al., 2023), aim to direct model attention towards 185
186 underrepresented classes, thereby improving model 186
187 performance in class-imbalanced settings. 187

3 Methodology

3.1 Problem Setup and Overview

We study low-resource and class-imbalanced transfer learning between heterogeneous clinical report datasets: a fully labeled source dataset \mathcal{D}_s and an unlabeled target dataset \mathcal{U}_t . Although the two datasets (domains) share some similar clinical knowledge, distribution shift and differences in label distribution make direct transfer challenging.

We formulate cross-domain adaptation as a budgeted active learning problem: given an annotation budget $B \ll N_t$, our goal is to select a small but high-utility subset $\mathcal{Q} \subset \mathcal{U}_t$. The selected samples are then annotated and merged with \mathcal{D}_s to form an expanded training set. With supervised fine-tuning on the final dataset, the knowledge can be effectively transferred and the model performance across both domains also improved.

The overall framework of RADS is shown in Figure 2. Our approach consists of three stages: (1) we train an active learner on \mathcal{D}_s and compute informativeness signals for \mathcal{U}_t via Monte-Carlo (MC) dropout; (2) we define a prior-aware utility that combines BALD-based mutual information (Houlsby et al., 2011) with pseudo-label class weighting to explicitly control the quality of selected samples for transfer learning under severe class imbalance; and (3) we train a reinforcement learning sampler to select samples that maximize prior-aware utility while discouraging redundant selections. The pseudocode for this part is provided in Appendix A.

3.2 Active Learner

We first fine-tune a lightweight classifier f_ϕ on the labeled source dataset \mathcal{D}_s . For each unlabeled target report in training pool $x \in \mathcal{U}_t$, we estimate epistemic uncertainty via MC dropout (Gal and Ghahramani, 2016). Specifically, we keep dropout activated at inference time and perform K stochastic forward passes. Each pass corresponds to sampling a dropout mask, yielding a sampled set of network weights \mathbf{w}_k and a predictive distribution:

$$p_k(y | x) = \text{softmax}(f_\phi(x; \mathbf{w}_k)) \quad (1)$$

Aggregating these K stochastic predictions approximates the posterior predictive distribution. We compute the MC predictive mean as:

$$\bar{p}(y | x) = \frac{1}{K} \sum_{k=1}^K p_k(y | x) \quad (2)$$

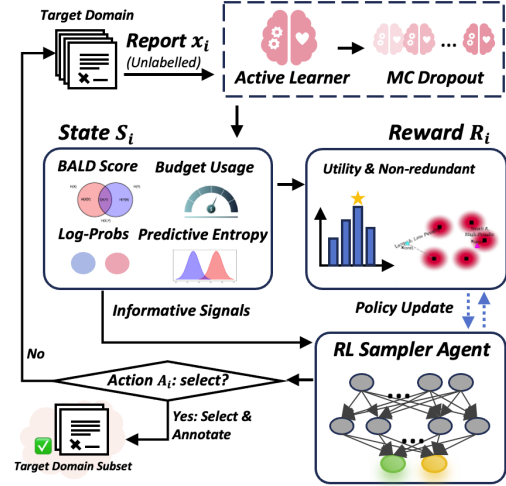


Figure 2: RADS framework for RL-based active sampling under domain shift. The active learner is fine-tuned on the source domain, and MC dropout is used to score unlabeled target reports and construct informativeness signals (state). An RL sampler then selects a subset for annotation by maximizing the reward, producing a target set for joint fine-tuning with the source data.

In addition, we retain the mean log-probability vector $\bar{\ell}(x) = \log \bar{p}(\cdot | x)$, which serves as a representation for redundancy estimation in our RL-based sampler (Section 3.5).

Based on this active learner, we also define a pseudo label $\hat{y}(x) = \arg \max_y \bar{p}(y | x)$ and estimate the predicted target class prior:

$$\hat{\pi}_+ = \frac{1}{N_t} \sum_{x \in \mathcal{U}_t} \mathbb{1}[\hat{y}(x) = 1], \quad (3)$$

$$\hat{\pi}_- = 1 - \hat{\pi}_+.$$

These priors let us correct selection bias when the pool is imbalanced or the source-trained model is miscalibrated on the target domain.

3.3 BALD Signal

To score informativeness for unlabeled target-domain samples, we use BALD, which quantifies the mutual information between the predicted label and the model parameters. Let $H(\cdot)$ denote entropy. For each $x \in \mathcal{U}_t$, we compute:

$$\text{PE}(x) = H(\bar{p}(\cdot | x)), \quad (4)$$

$$\text{EE}(x) = \frac{1}{K} \sum_{k=1}^K H(p_k(\cdot | x)), \quad (5)$$

$$\text{MI}(x) = \text{PE}(x) - \text{EE}(x). \quad (6)$$

Here, $\text{MI}(x)$ is the BALD score. It is large when the predictive distribution is uncertain overall (high

PE) while individual stochastic models are relatively confident but disagree with each other (low EE). We normalize $\text{MI}(x)$ to $[0, 1]$ over \mathcal{U}_t , denoted as $\widetilde{\text{MI}}(x)$.

We treat samples with high $\widetilde{\text{MI}}(x)$ as informative and assign them higher utility in our selection policy. Prioritizing these samples for annotation is expected to reduce the model’s uncertainty and improve transfer to the target domain.

3.4 Prior-Aware Utility for Sample Selection

Selecting the top- B uncertain samples can sometimes produce an extreme class skew. This often happens under domain shift and severe class imbalance, where the source-trained active learner may predict biased pseudo labels on the target domain. To control the selected class mixture, we introduce a prior-aware utility. We define class weights using the estimated prior:

$$\begin{aligned} w_+ &= \frac{\rho}{\text{clip}(\hat{\pi}_+)}, \\ w_- &= \frac{1 - \rho}{1 - \text{clip}(\hat{\pi}_+)}. \end{aligned} \quad (7)$$

where $\text{clip}(\cdot)$ clamps probabilities away from $\{0, 1\}$ for stability and ρ is a hyperparameter that trades off class-balance control and informativeness. We then define the utility:

$$u(x) = \widetilde{\text{MI}}(x) \cdot \begin{cases} w_+, & \hat{y}(x) = 1, \\ w_-, & \hat{y}(x) = 0. \end{cases} \quad (8)$$

This utility favors informative samples and shifts selection toward the desired class ratio.

3.5 RL-based Sample Selection Strategy

At each step t , the sampler agent observes the current candidate x_t and decides whether to select or discard. An episode ends when B samples are selected or the pool is exhausted.

State. For each candidate x_t , the state vector combines the active learner signals and a budget progress term:

$$s_t = \left[\bar{\ell}(x_t); \text{PE}(x_t); \text{MI}(x_t); |S_t|/B \right] \quad (9)$$

where $\bar{\ell}(x_t)$ is the mean log-probability vector computed from MC dropout; $\text{PE}(x_t)$ is the predictive entropy; $\text{MI}(x_t)$ is the BALD score; and $|S_t|/B$ indicates the fraction of the annotation budget already consumed, with S_t denoting the set of selected samples so far. In our binary setting, $\bar{\ell}(x_t) \in \mathbb{R}^2$, hence the overall dimension of the state vector is 5.

Reward. Our reward encourages the agent to select samples that are both (i) informative for learning under class imbalance and (ii) non-redundant with respect to previously selected instances. Specifically, when the agent selects the current candidate ($a_t = 1$) and the budget is not yet exhausted ($|S_t| < B$), we define:

$$r_t = u(x_t) - \lambda \cdot \text{Red}(x_t, S_t) \quad (10)$$

and set $r_t = 0$ otherwise. Here, $u(x_t)$ is the prior-aware utility (Section 3.4) and λ controls the strength of the diversity regularization.

To discourage selecting near-duplicate samples, we measure redundancy in the active learner’s predictive representation space. For a candidate x and the current selected set S , we first compute the distance to its nearest selected neighbor:

$$\delta(x, S) = \begin{cases} +\infty, & |S| = 0, \\ \min_{x' \in S} \|\bar{\ell}(x) - \bar{\ell}(x')\|_2, & \text{otherwise.} \end{cases} \quad (11)$$

We then convert this distance into a bounded redundancy score:

$$\text{Red}(x, S) = \begin{cases} 0, & |S| = 0, \\ \frac{1}{1 + \delta(x, S)}, & \text{otherwise.} \end{cases} \quad (12)$$

This definition yields a larger penalty when x is very close to an existing selection (small δ), and a smaller penalty when x is far away (large δ). As a result, the agent is encouraged to select diverse samples while still prioritizing high-utility ones.

Dueling DQN Sampler Agent. We learn a Q-function $Q_\theta(s, a)$ with a dueling DQN architecture (Wang et al., 2016) and optimize it via the standard DQN objective (Mnih et al., 2015). We maintain an experience replay buffer \mathcal{B} and a target network Q_{θ^-} . At each gradient step, we minimize the temporal-difference loss:

$$\begin{aligned} \mathcal{L}(\theta) &= \mathbb{E}_{(s,a,r,s',d) \sim \mathcal{B}} \left[(Q_\theta(s, a) - y)^2 \right], \\ y &= r + \gamma(1 - d) \max_{a'} Q_{\theta^-}(s', a'). \end{aligned} \quad (13)$$

where γ is the discount factor and $d \in \{0, 1\}$ indicates episode termination. We adopt ϵ -greedy exploration with a decaying ϵ schedule, periodically synchronize θ^- with θ , and finally use the learned policy $\pi(s) = \arg \max_a Q_\theta(s, a)$ to select B samples from \mathcal{U}_t .

Attribute	CHIFIR	PIFIR	MIMIC-CXR (subset)
Report Type	Cytology / Histopathology	PET-CT Radiology	Chest X-ray Radiology
Target Disease	Invasive Fungal Infection (IFI)	Invasive Fungal Infection (IFI)	Pneumonia
Label Type	Gold (manual)	Gold (manual)	Silver (auto-derived)
Class distribution (P/N)	14% / 86%	69% / 31%	39% / 61%
Dataset Size	283 reports (small)	201 reports (small)	493 reports (medium)

Table 1: Key attributes of the three datasets used in this study. P/N = positive/negative class proportions

4 Experimental Setup

4.1 Benchmark Datasets

We chose three real world clinical datasets (Rozova et al., 2023a, 2025; Johnson et al., 2019) as benchmarks in this study: the PET-CT Invasive Fungal Infection Reports corpus (PIFIR¹), the Cytology and Histopathology IFI Reports corpus (CHIFIR²), and the MIMIC Chest X-ray corpus (MIMIC-CXR³).

CHIFIR and PIFIR datasets are related to Invasive Fungal Infection (IFI), but the vocabulary used varies across them. The cytology and histopathology reports of the CHIFIR dataset assess tissue or fluid samples and describe the microscopic visualization of fungal organisms. The PET-CT reports from PIFIR assess metabolic activity and discuss the anatomical and morphological features of fungal lesions via PET imaging.

To assess transfer beyond IFI and beyond pathology-style reports, we also include MIMIC-CXR, a corpus of chest X-ray reports. We construct a *Pneumonia* subset by selecting the top 3,000 reports that are labeled as pneumonia by CheXpert’s (Irvin et al., 2019) weak labels. Although PIFIR and MIMIC-CXR both consist of radiology reports, they still differ greatly in reporting style and clinical phrasing. Moreover, pneumonia and IFI reflect distinct clinical contexts, further increasing the domain shift. Figure 3 shows differences in predominant clinical terms across the three datasets.



Figure 3: Word clouds for the CHIFIR (left), PIFIR (middle), and MIMIC-CXR (right) datasets. Word size corresponds to term frequency.

¹Available for credentialed users at <https://physionet.org/content/pifir/1.0.0/>

²Available for credentialed users at <https://physionet.org/content/corpus-fungal-infections/1.0.2/>

³Available for credentialed users at <https://physionet.org/content/mimic-cxr/2.1.0/>

All three datasets exhibit class imbalance. CHIFIR and MIMIC-CXR are dominated by negative cases, whereas PIFIR is dominated by positive cases. Throughout the paper we focus on transferring from other sources to PIFIR, and we provide results of other directions in the Appendix.

Table 1 summarizes the key characteristics of each dataset and highlights the challenges for transfer learning across them. Details of the dataset split are provided in Appendix B.

4.2 Evaluation Metrics

We evaluate performance using accuracy, F1 score, precision, recall, and ROC-AUC. Class imbalance in benchmark datasets makes the F1 score particularly important. Recall is also important, given that it is critical not to miss positive cases.

4.3 Baselines

We select the fine-tuned ClinicalBERT approach from previous work (Han et al., 2025) as the baseline. Table 2 shows the baseline results and reveals the challenges of knowledge transferability between these datasets. Models perform well when fine-tuned and evaluated on the same dataset. Without transfer learning, evaluation on a similar but still different dataset results in a clear performance drop. Although training on all datasets together can improve performance, it requires annotating all reports, which is labor-intensive. Full reproducibility details can be found in Appendix B.

5 Experimental Results

5.1 Transfer Learning Performance

We compare our method with several other active learning approaches to analyze the impact of different sample selection methods on knowledge transfer performance. 1) **Random Selection**: Randomly select samples from the unlabeled target domain. Each experiment is run five times to reduce variance and obtain more reliable results. We report the mean evaluation metrics over these five runs. 2) **Uncertainty-based Selection** (Nguyen

Transfer Learning to PIFIR Strategy		Performance on PIFIR				Performance on CHIFIR				Performance on MIMIC-CXR			
Datasets		Acc	F1	P	R	Acc	F1	P	R	Acc	F1	P	R
Baseline	PIFIR	0.714	0.812	0.788	0.839	–	–	–	–	–	–	–	–
Zero-shot	CHIFIR	0.357	0.229	1.000	0.129	0.942	0.824	0.778	0.875	–	–	–	–
	MIMIC-CXR	0.738	0.841	0.763	0.935	–	–	–	–	0.859	0.811	0.833	0.789
Full-shot	CHIFIR + PIFIR	0.857	0.900	0.931	0.871	0.904	0.615	0.800	0.500	–	–	–	–
	MIMIC-CXR + PIFIR	0.833	0.889	0.875	0.903	–	–	–	–	0.848	0.800	0.811	0.789

Table 2: Performance comparison of ClinicalBERT under different transfer strategies. Zero-shot transfer refers to fine-tuning solely on the source dataset. We set this as our baseline. Full-shot transfer refers to jointly fine-tuning on both source and target datasets. We assume this represents the best possible transfer learning performance.

et al., 2022): Selects k samples by predictive uncertainty (lowest confidence) from the active learner. 3) **Diversity-based Selection** (Yuan et al., 2020): Selects the k most diverse samples by calculating the cosine distance between each report embedding in the unlabeled dataset and the embeddings in the labeled dataset. 4) **LM-DPP Selection** (Wang et al., 2024): This method jointly models uncertainty and diversity using a Determinantal Point Process (DPP) kernel. Following the original work, we set the trade-off coefficient between uncertainty and diversity to 0.5 and select the subset of size k that maximizes the DPP objective for annotation. 5) **TAGCOS Selection** (Zhang et al., 2025): A task-agnostic selection baseline that selects k samples according to its gradient-based selection criterion. 6) **BatchBALD Selection** (Kirsch et al., 2019): Selects k samples using a batch acquisition strategy that extends BALD by maximizing joint mutual information under MC dropout.

Table 3 reports transfer learning results from CHIFIR to PIFIR. Uncertainty-, diversity-, and LM-DPP-based selection yield comparable or worse performance than random sampling. While TAGCOS and BatchBALD attain relatively high F1 scores on PIFIR, their ROC-AUC is noticeably lower. In contrast, our method RADS achieves the best performance on PIFIR while maintaining competitive performance on the source domain (CHIFIR). This indicates strong sample efficiency, requiring only $5/135 \approx 3.7\%$ of the target training set to obtain substantial transfer gains.

Table 4 reports transfer learning results from MIMIC-CXR to PIFIR. Other baselines provide limited gains and are often on par with or below random selection. RADS achieves better target performance (F1 on PIFIR = 0.882) and remains highly sample-efficient, requiring annotation of only $2/135 \approx 1.5\%$ of the target training set.

We also conduct transfer learning experiments

from PIFIR to CHIFIR. The results show that our method still achieves the best performance with only 8 samples selected from target dataset CHIFIR. More detailed discussion is shown in Appendix D.

RADS consistently outperforms strong baselines, demonstrating superior sample-efficient transfer performance. More robustness analysis under imbalanced settings appears in Appendix C.

5.2 Learning Curves under Varying Budgets

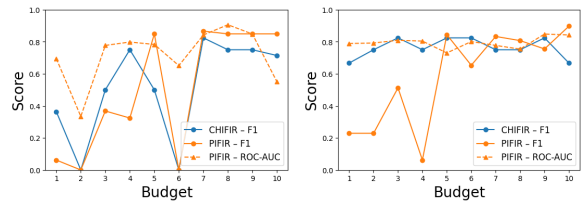


Figure 4: Transfer from CHIFIR to PIFIR under baselines BatchBALD (left) and TAGCOS (right).

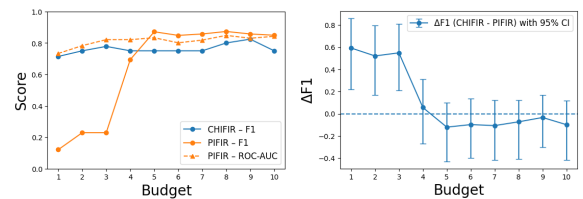


Figure 5: Transfer from CHIFIR to PIFIR under our method RADS.

We analyze the effect of annotation budget on transfer performance. Figure 4 shows the transfer performance from CHIFIR to PIFIR across budgets under two strong baselines. BatchBALD is highly unstable; a small change in budget can flip the model from almost perfect to completely broken. TAGCOS is more stable but still remains unreliable at small budgets. Figure 5 shows the transfer performance from CHIFIR to PIFIR with our method. The left graph shows that starting from budget 5, F1 on PIFIR stays around 0.85–0.87 and additional

Knowledge Transfer from CHIFIR to PIFIR												
Strategy	Performance on PIFIR					Performance on CHIFIR					Transfer Gap	
	Accuracy	F1-score	Precision	Recall	ROC-AUC	Accuracy	F1-score	Precision	Recall	ROC-AUC	$\Delta F1$	95% CI
Random	0.595	0.639	0.885	0.561	0.813	0.927	0.746	0.805	0.700	0.938	–	–
Uncertainty	0.524	0.545	0.923	0.387	0.830	0.942	0.824	0.778	0.875	0.977	0.278	[-0.037, 0.530]
Diversity	0.595	0.638	0.938	0.484	0.809	0.942	0.800	0.857	0.750	0.974	0.162	[-0.167, 0.412]
LM-DPP	0.571	0.609	0.933	0.452	0.839	0.904	0.615	0.800	0.500	0.977	0.007	[-0.418, 0.319]
TAGCOS	0.762	0.844	0.818	0.871	0.730	0.942	0.824	0.778	0.875	0.972	-0.020	[-0.310, 0.168]
BatchBALD	0.738	0.849	0.738	1.000	0.783	0.885	0.500	0.750	0.375	0.946	-0.349	[-0.833, -0.019]
RADS	0.810	0.871	0.871	0.871	0.833	0.923	0.750	0.750	0.750	0.977	-0.121	[-0.430, 0.100]

Table 3: Transfer learning performance from CHIFIR to PIFIR with 5 samples selected in PIFIR under different sample selection strategies. $\Delta F1 = F1(\text{CHIFIR}) - F1(\text{PIFIR})$. CI = Confidence Interval.

Knowledge Transfer from MIMIC-CXR to PIFIR												
Strategy	Performance on PIFIR					Performance on MIMIC-CXR					Transfer Gap	
	Accuracy	F1-score	Precision	Recall	ROC-AUC	Accuracy	F1-score	Precision	Recall	ROC-AUC	$\Delta F1$	95% CI
Random	0.766	0.848	0.808	0.842	0.806	0.862	0.824	0.808	0.842	0.932	–	–
Uncertainty	0.738	0.831	0.794	0.871	0.636	0.848	0.810	0.780	0.842	0.916	-0.021	[-0.165, 0.126]
Diversity	0.738	0.845	0.750	0.968	0.616	0.859	0.816	0.816	0.816	0.933	-0.029	[-0.157, 0.102]
LM-DPP	0.738	0.831	0.794	0.871	0.636	0.848	0.810	0.780	0.842	0.916	-0.021	[-0.165, 0.126]
TAGCOS	0.738	0.836	0.778	0.903	0.795	0.862	0.831	0.821	0.842	0.930	-0.005	[-0.131, 0.139]
BatchBALD	0.762	0.848	0.800	0.903	0.827	0.889	0.861	0.829	0.895	0.949	0.012	[-0.107, 0.141]
RADS	0.810	0.882	0.811	0.968	0.880	0.869	0.840	0.791	0.895	0.921	-0.043	[-0.153, 0.072]

Table 4: Transfer learning performance from MIMIC-CXR to PIFIR with 2 samples selected in PIFIR under different sample selection strategies. $\Delta F1 = F1(\text{MIMIC-CXR}) - F1(\text{PIFIR})$. CI = Confidence Interval.

labels provide marginal improvements, while CHIFIR performance remains stable. The right graph plots the domain gap $\Delta F1$ versus budget with 95% confidence intervals. From budget = 4 onward, the gap is effectively closed ($\Delta F1 \approx 0$ with overlapping confidence intervals), suggesting the selected subset suffices to eliminate the transfer gap.

As MIMIC-CXR is larger than CHIFIR, identifying informative target samples is less challenging for all models, even under low annotation budgets. The transfer performance from MIMIC-CXR to PIFIR across budgets is shown in Appendix E.

5.3 Ablation Study

To evaluate the effectiveness of each component in RADS, we conduct ablation studies as shown in Table 5. Replacing the RL sampler with a greedy selector (No RL) leads to a clear drop performance. Although this variant optimizes the same objective, it lacks the sequential decision-making needed to balance exploration and redundancy control. Selecting samples by BALD signal (MI Only) fails, implying that uncertainty-only criteria can favor noisy or out-of-distribution target examples under domain shift. Using the prior-aware utility (Utility Only) greatly improves results, confirming the benefit of class- and quality-aware selection, but remains below our method, highlighting the additional gains from discouraging redundant selec-

tions. Finally, RADS further improves over Utility Only, demonstrating that the RL sampler provides additional gains by learning a non-redundant, globally optimized subset selection policy rather than relying on pointwise ranking.

Method	Accuracy	F1-score	Precision	Recall	ROC-AUC
No RL	0.571	0.609	0.933	0.452	0.798
MI Only	0.262	0.000	0.000	0.000	0.475
Utility Only	0.786	0.866	0.806	0.935	0.833
RADS	0.810	0.871	0.871	0.871	0.833

Table 5: Ablation results under the hard transfer setting from CHIFIR to PIFIR with a labeling budget of 5.

5.4 Selected Sample Quality Analysis

We audit the quality of the selected target samples. Figure 6 shows that in CHIFIR to PIFIR transfer, the source-trained pseudo labels are notably misaligned with the annotated labels, yet RADS still selects mostly true positives, which helps improve target performance. In MIMIC-CXR to PIFIR, the pseudo-label ratio better matches the annotated composition, consistent with smaller domain shift and improved calibration on the target domain.

Figure 7 visualizes the CHIFIR to PIFIR (left) and MIMIC-CXR to PIFIR (right) transfer, where our method selects 5 samples for adaptation. Before transfer learning, the decision boundary only captures the source dataset specific separation. Af-

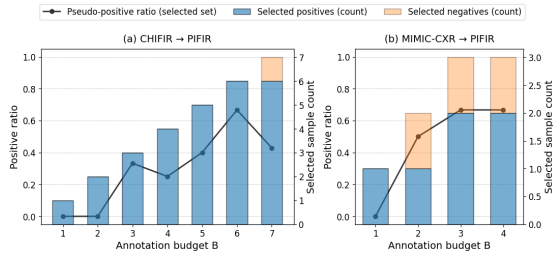


Figure 6: Selected sample analysis under CHIFIR to PIFIR (left) and MIMIC-CXR to PIFIR (right) transfer. The black line shows the pseudo-positive ratio predicted by the source-trained active learner, and the bars report the numbers of true positives and true negatives after manual annotation (blue/yellow).

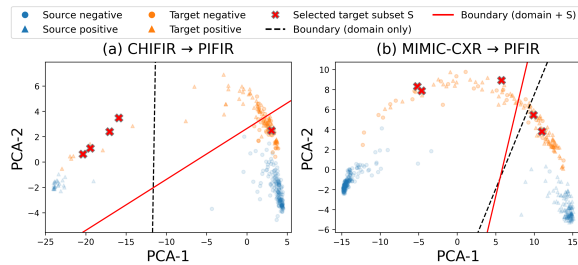


Figure 7: Transfer learning from CHIFIR to PIFIR (left) and MIMIC-CXR to PIFIR (right) with PCA projection of report embeddings. Red \times markers denote the selected PIFIR subset S . The dashed black line shows the decision boundary learned from the source dataset only. The solid red line shows the boundary after augmenting with selected samples. Source = CHIFIR / MIMIC-CXR dataset, Target = PIFIR dataset.

495 ter adding the selected target subset, the boundary
 496 rotates and shifts toward a direction that better re-
 497 flects the class layout from both datasets.

498 5.5 Efficiency and Budget-Aware Transfer

499 **Runtime Efficiency** Our RL-based sampler in
 500 RADS is light-weight, requiring only a few sec-
 501 onds to produce a selection. Selecting 2 samples
 502 (MIMIC-CXR to PIFIR transfer) takes around 3
 503 seconds, and selecting 5 samples (CHIFIR to PI-
 504 FIR transfer) takes around 9 seconds. Given that
 505 annotating a single report takes about one minute,
 506 the selection overhead is negligible. Compared
 507 with full-shot transfer that labels the entire target
 508 training set (135 reports), RADS achieves compa-
 509 rable target performance with only 2 or 5 annotated
 510 reports.

511 **Budget Utilization and Early Stopping** As
 512 our RL-based sampler encodes budget progress
 513 ($|S_t|/B$) in the state, it can also provide guidance
 514 on how many samples are worth annotating during

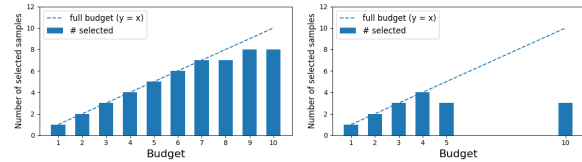


Figure 8: Number of selected PIFIR samples versus the annotation budget B when training on CHIFIR to PIFIR (left) and MIMIC-CXR to PIFIR (right).

515 transfer. Figure 8 shows the actual number of PI-
 516 FIR samples selected as the budget increases. For
 517 CHIFIR to PIFIR transfer, the policy stops fully
 518 consuming the budget once $B \geq 8$, consistent with
 519 our results showing good transfer performance al-
 520 ready at $B = 5$. For MIMIC-CXR to PIFIR trans-
 521 fer, the policy no longer uses the full budget when
 522 $B \geq 5$, aligning with good target performance
 523 achieved with $B = 2$ labeled PIFIR samples. This
 524 pattern is useful as in active learning it is important
 525 to know when it is time to stop adding samples.

526 5.6 Transfer Gap between Datasets

527 To explain why model transfer from MIMIC-CXR
 528 to PIFIR is easier than CHIFIR to PIFIR, we further
 529 analyze the distribution gap between these datasets.

530 We quantify overlap between datasets with a
 531 shared unigram–bigram vocabulary (Elangovan
 532 et al., 2021). Coverage of PIFIR-test n-grams is
 533 higher for MIMIC-CXR to PIFIR than for CHIFIR
 534 to PIFIR (0.193 vs. 0.115). The Jaccard similar-
 535 ity between source and target vocabularies is also
 536 higher for MIMIC-CXR to PIFIR than for CHIFIR
 537 to PIFIR (0.187 vs. 0.124). This suggests a smaller
 538 lexical domain shift between MIMIC-CXR and PI-
 539 FIR, as it is illustrated by our empirical results fine-
 540 tuning MIMIC-CXR. More detailed differences are
 541 analyzed in Appendix F.

542 6 Conclusion

543 In this work, we studied transfer learning for
 544 disease detection under low-resource and class-
 545 imbalanced conditions. We proposed RADS, an
 546 RL-based sampler that jointly optimizes a prior-
 547 aware utility for class-mixture control and a diver-
 548 sity regularizer to avoid near-duplicate selections.
 549 Our approach improves model performance and
 550 adaptability across medical datasets compared to
 551 traditional sample selection strategies. We expect
 552 this approach to generalize to other transfer learn-
 553 ing problems in clinical NLP, offering a promising
 554 direction for broader validation and impact.

555 Limitations

556 Despite demonstrating promising results, our ap-
557 proach has several limitations. The effectiveness of
558 our RL-based sample selection heavily depends on
559 the feedback provided by the active learner. This
560 places high quality demands on the original gold
561 dataset. Our formulation controls class mixture
562 only through predicted priors and does not explic-
563 itly incorporate richer clinical knowledge. Future
564 work will study more robust uncertainty estima-
565 tion, and investigate methods to integrate domain
566 knowledge to further improve reliability and gener-
567 alization.

568 Acknowledgments

569 Supported by anonymous grant and relevant Ethics
570 approvals.

571 References

572 Naif Radi Aljohani, Ayman Fayoumi, and Saeed-Ul Has-
573 san. 2023. A novel focal-loss and class-weight-aware
574 convolutional neural network for the classification
575 of in-text citations. *Journal of Information Science*,
576 49(1):79–92.

577 Zaid Alyafeai, Maged Saeed AlShaibani, and Irfan
578 Ahmad. 2020. A survey on transfer learning
579 in natural language processing. *arXiv preprint*
580 *arXiv:2007.04239*.

581 Imane Araf, Ali Idri, and Ikram Chairi. 2024. Cost-
582 sensitive learning for imbalanced medical data: a
583 review. *Artificial Intelligence Review*, 57(4):80.

584 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
585 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
586 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
587 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
588 Gretchen Krueger, Tom Henighan, Rewon Child,
589 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
590 Winter, and 12 others. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

594 Sergio Consoli, Peter Markov, Nikolaos I Stilianakis,
595 Lorenzo Bertolini, Antonio Puertas Gallardo, and
596 Mario Ceresa. 2024. Epidemic information extrac-
597 tion for event-based surveillance using large lan-
598 guage models. In *International Congress on Informa-
599 tion and Communication Technology*, pages 241–252. Springer Nature Singapore Singapore.

601 Ricardo C Cury, Istvan Megyeri, Tony Lindsey, Rob-
602 son Macedo, Juan Batlle, Shwan Kim, Brian Baker,
603 Robert Harris, and Reese H Clark. 2021. Natural
604 Language Processing and Machine Learning for De-
605 tecting of Respiratory Illness by Chest CT Imag-
606 ing and Tracking of COVID-19 Pandemic in the

United States. *Radiology: Cardiothoracic Imaging*,
3(1):e200596.

Aparna Elangovan, Jiayuan He, and Karin Verspoor. 2021. [Memorization vs. generalization : Quantifying data leakage in NLP performance evaluation](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1325–1335, Online. Association for Computational Linguistics.

Meng Fang, Yuan Li, and Trevor Cohn. 2017. [Learning how to active learn: A deep reinforcement learning approach](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 595–605, Copenhagen, Denmark. Association for Computational Linguistics.

Yifan Fu, Xingquan Zhu, and Bin Li. 2013. A survey on instance selection for active learning. *Knowledge and information systems*, 35(2):249–283.

Yarin Gal and Zoubin Ghahramani. 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR.

Kushankur Ghosh, Colin Bellinger, Roberto Corizzo, Paula Branco, Bartosz Krawczyk, and Nathalie Japkowicz. 2024. The class imbalance problem in deep learning. *Machine Learning*, 113(7):4845–4901.

Julius Gonsior, Christian Falkenberg, Silvio Magino, Anja Reusch, Claudio Hartmann, Maik Thiele, and Wolfgang Lehner. 2024. Comparing and improving active learning uncertainty measures for transformer models by discarding outliers. *Information systems frontiers*, pages 1–17.

Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang. 2022. [PPT: Pre-trained prompt tuning for few-shot learning](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8410–8423, Dublin, Ireland. Association for Computational Linguistics.

Hairani Hairani, Triyanna Widiyaningtyas, and Didik Dwi Prasetya. 2024. Addressing class imbalance of health data: A systematic literature review on modified synthetic minority oversampling technique (smote) strategies. *JOIV: International Journal on Informatics Visualization*, 8(3):1310–1318.

Wei Han, David Martinez, Vlada Rozova, Lawrence Cavedon, Anna Khanina, Leon J Worth, Monica A Slavin, Karin A Thursky, and Karin Verspoor. 2025. Automated detection of invasive fungal infections in clinical reports using medical language models. *Studies in Health Technology and Informatics*, 329:1002–1007.

Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. 2011. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*.

662	Kexin Huang, Jaan Altsosaar, and Rajesh Ranganath.	Mieradilijiang Maimaiti, Yang Liu, Huanbo Luan, and	717
663	2019. ClinicalBERT: Modeling clinical notes and	Maosong Sun. 2021. Enriching the transfer learning	718
664	predicting hospital readmission. <i>CHIL 2020 Work-</i>	with pre-trained lexicon embedding for low-resource	719
665	<i>shop</i> .	neural machine translation. <i>Tsinghua Science and</i>	720
		<i>Technology</i> , 27(1):150–163.	721
666	Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu,	David Martinez, Michelle R Ananda-Rajah, Hanna	722
667	Silviana Ciurea-Ilcus, Chris Chute, Henrik Mark-	Suominen, Monica A Slavin, Karin A Thursky, and	723
668	lund, Behzad Haghgoo, Robyn Ball, Katie Shpan-	Lawrence Cavedon. 2015. Automatic detection of pa-	724
669	skaya, and 1 others. 2019. Chexpert: A large chest	tients with invasive fungal disease from free-text com-	725
670	radiograph dataset with uncertainty labels and expert	puted tomography (CT) scans. <i>Journal of Biomedical</i>	726
671	comparison. In <i>Proceedings of the AAAI conference</i>	<i>Informatics</i> , 53:251–260.	727
672	<i>on artificial intelligence</i> , volume 33, pages 590–597.		
673	Henrik Elvang Jensen. 2021. Histopathology in the di-	Volodymyr Mnih, Koray Kavukcuoglu, David Silver,	728
674	agnosis of invasive fungal diseases. <i>Current Fungal</i>	Andrei A Rusu, Joel Veness, Marc G Bellemare,	729
675	<i>Infection Reports</i> , 15(1):23–31.	Alex Graves, Martin Riedmiller, Andreas K Fidje-	730
		land, Georg Ostrovski, and 1 others. 2015. Human-	731
676	Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz,	level control through deep reinforcement learning.	732
677	Nathaniel R Greenbaum, Matthew P Lungren, Chih-	<i>Nature</i> , 518(7540):529–533.	733
678	ying Deng, Roger G Mark, and Steven Horng.		
679	2019. MIMIC-CXR, a de-identified publicly avail-	Vu-Linh Nguyen, Mohammad Hossein Shaker, and	734
680	able database of chest radiographs with free-text re-	Eyke Hüllermeier. 2022. How to measure uncertainty	735
681	ports. <i>Scientific data</i> , 6(1):317.	in uncertainty sampling for active learning. <i>Machine</i>	736
		<i>Learning</i> , 111(1):89–122.	737
682	Justin M Johnson and Taghi M Khoshgoftaar. 2019. Sur-	V. Rozova, A. Khanina, J. Ong, R. Alipour, L. Worth,	738
683	vey on deep learning with class imbalance. <i>Journal</i>	M. Slavin, K. Thursky, and K. Verspoor. 2025. PIFIR:	739
684	<i>of Big Data</i> , 6(1):1–54.	PET-CT Invasive Fungal Infection Reports (version	740
		1.0.0) . [Data set].	741
685	Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal.	V. Rozova, A. Khanina, J. Teng, J. Teh, L. Worth,	742
686	2019. Batchbald: Efficient and diverse batch acqui-	M. Slavin, K. Thursky, and K. Verspoor. 2023a. CHI-	743
687	sition for deep bayesian active learning. <i>Advances in</i>	FIR: Cytology and Histopathology Invasive Fungal	744
688	<i>neural information processing systems</i> , 32.	Infection Reports (version 1.0.0) . [Data set].	745
689	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon	Vlada Rozova, Anna Khanina, Jasmine C Teng,	746
690	Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang.	Joanne SK Teh, Leon J Worth, Monica A Slavin,	747
691	2020. BioBERT: a pre-trained biomedical language	Karin A Thursky, and Karin Verspoor. 2023b. Detect-	748
692	representation model for biomedical text mining.	ing evidence of invasive fungal infections in cytology	749
693	<i>Bioinformatics</i> , 36(4):1234–1240.	and histopathology reports enriched with concept-	750
694	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	level annotations. <i>Journal of Biomedical Informatics</i> ,	751
695	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	139:104293.	752
696	Deng, Chenyu Zhang, Chong Ruan, and 1 others.		
697	2024a. DeepSeek-V3 Technical Report. <i>arXiv</i>	Chuanqi Tan, Fuchun Sun, Tao Kong, Wenchang Zhang,	753
698	<i>preprint arXiv:2412.19437</i> .	Chao Yang, and Chunfang Liu. 2018. A survey on	754
		deep transfer learning. In <i>Artificial Neural Networks</i>	755
699	Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mo-	<i>and Machine Learning–ICANN 2018: 27th Inter-</i>	756
700	hta, Tenghao Huang, Mohit Bansal, and Colin A Raf-	<i>national Conference on Artificial Neural Networks,</i>	757
701	ffel. 2022. Few-shot parameter-efficient fine-tuning	<i>Rhodes, Greece, October 4-7, 2018, Proceedings,</i>	758
702	is better and cheaper than in-context learning. <i>Ad-</i>	<i>Part III 27</i> , pages 270–279. Springer.	759
703	<i>vances in Neural Information Processing Systems</i> ,		
704	35:1950–1965.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	760
		Martinet, Marie-Anne Lachaux, Timothée Lacroix,	761
705	Ying Liu, Haozhu Wang, Huixue Zhou, Mingchen Li,	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	762
706	Yu Hou, Sicheng Zhou, Fang Wang, Rama Hoetzlein,	Azhar, and 1 others. 2023. Llama: Open and effi-	763
707	and Rui Zhang. 2024b. A review of reinforcement	cient foundation language models. <i>arXiv preprint</i>	764
708	learning for natural language processing and applica-	<i>arXiv:2302.13971</i> .	765
709	tions in healthcare. <i>Journal of the American Medical</i>		
710	<i>Informatics Association</i> , 31(10):2379–2393.	David W Townsend, Jonathan PJ Carney, Jeffrey T Yap,	766
		and Nathan C Hall. 2004. PET/CT today and tomor-	767
711	Pilar López-Úbeda, Manuel Carlos Díaz-Galiano,	row. <i>Journal of Nuclear Medicine</i> , 45(1 suppl):4S–	768
712	Teodoro Martín-Noguerol, Antonio Luna, L Alfonso	14S.	769
713	Ureña-López, and M Teresa Martín-Valdivia. 2020.		
714	COVID-19 detection in radiological text reports in-	Peng Wang, Xiaobin Wang, Chao Lou, Shengyu Mao,	770
715	tegrating entity recognition. <i>Computers in Biology</i>	Pengjun Xie, and Yong Jiang. 2024. Effective demon-	771
716	<i>and Medicine</i> , 127:104066.	stration annotation for in-context learning via lan-	772
		guage model-based determinantal point process . In	773

774 *Proceedings of the 2024 Conference on Empirical*
775 *Methods in Natural Language Processing*, pages
776 1266–1280, Miami, Florida, USA. Association for
777 Computational Linguistics.

778 Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt,
779 Marc Lanctot, and Nando Freitas. 2016. Dueling
780 network architectures for deep reinforcement learn-
781 ing. In *International conference on machine learning*,
782 pages 1995–2003. PMLR.

783 Karl Weiss, Taghi M Khoshgoftaar, and DingDing
784 Wang. 2016. A survey of transfer learning. *Jour-
785 nal of Big Data*, 3:1–40.

786 Cynthia Yang, Egill A Fridgeirsson, Jan A Kors,
787 Jenna M Reps, and Peter R Rijnbeek. 2024. Impact
788 of random oversampling and random undersampling
789 on the performance of prediction models developed
790 using observational health data. *Journal of Big Data*,
791 11(1):7.

792 Yi Yang, Zhigang Ma, Feiping Nie, Xiaojun Chang, and
793 Alexander G Hauptmann. 2015. Multi-class active
794 learning by uncertainty sampling with diversity max-
795 imization. *International Journal of Computer Vision*,
796 113:113–127.

797 Michelle Yuan, Hsuan-Tien Lin, and Jordan Boyd-
798 Graber. 2020. Cold-start active learning through self-
799 supervised language modeling. In *Proceedings of the*
800 *2020 Conference on Empirical Methods in Natural*
801 *Language Processing (EMNLP)*, pages 7935–7948,
802 Online. Association for Computational Linguistics.

803 Jipeng Zhang, Yaxuan Qin, Renjie Pi, Weizhong Zhang,
804 Rui Pan, and Tong Zhang. 2025. TAGCOS: Task-
805 agnostic gradient clustered coreset selection for in-
806 struction tuning data. In *Findings of the Association*
807 *for Computational Linguistics: NAACL 2025*, pages
808 4671–4686.

A Algorithm Pseudocode

We provide the pseudocode for the strategy of sam-
ple selection in our approach below.

Require: pool \mathcal{U}_t , budget B , episodes N , utility $u(\cdot)$, diver-
sity weight λ
feature set $\mathcal{F} = \{\bar{\ell}(x), \text{PE}(x), \widetilde{\text{MI}}(x), |S|/B\}$, ac-
tion set $\mathcal{A} = \{0, 1\}$
1: $\text{env} \leftarrow \text{RLSAMPLESELECTIONENV}(\mathcal{U}_t, \mathcal{F}, u, B, \lambda)$
2: Initialize online network Q_ϕ and target network $\hat{Q}_\phi \leftarrow$
 Q_ϕ
3: Initialize replay buffer $\mathcal{D} \leftarrow \emptyset$
4: Initialize exploration rate ϵ
5: **for** $\text{episode} = 1$ **to** N **do**
6: $s \leftarrow \text{env.reset}()$; $\text{done} \leftarrow \text{false}$
7: **while not** done **do**
8: $a \leftarrow \text{EPSGREEDY}(Q_\phi, s, \epsilon)$
9: $(s', r, \text{done}) \leftarrow \text{env.step}(a)$
10: Add $(s, a, r, s', \text{done})$ to \mathcal{D}
11: **if** $|\mathcal{D}| \geq M$ **then** $\triangleright M$ is minibatch size
12: $\mathcal{M} \leftarrow \text{SAMPLEMINIBATCH}(\mathcal{D}, M)$
13: $\text{UPDATENETS}(Q_\phi, \hat{Q}_\phi, \mathcal{M})$
14: **end if**
15: $s \leftarrow s'$
16: **end while**
17: **if** $\text{episode} \bmod K_{\text{upd}} = 0$ **then**
18: $\hat{Q}_\phi \leftarrow Q_\phi$
19: **end if**
20: $\epsilon \leftarrow \text{DECAYEPS}(\epsilon)$
21: **end for**
Selection (Greedy Policy)
22: $\mathcal{S} \leftarrow \emptyset$; $s \leftarrow \text{env.reset}()$; $\text{done} \leftarrow \text{false}$
23: **while not** done **do**
24: $\text{id} \leftarrow \text{env.currentId}()$
25: $a \leftarrow \arg \max_{a' \in \mathcal{A}} Q_\phi(s, a')$
26: $(s', _, \text{done}) \leftarrow \text{env.step}(a)$
27: **if** $a = 1$ **then**
28: $\mathcal{S} \leftarrow \mathcal{S} \cup \{\text{id}\}$
29: **end if**
30: $s \leftarrow s'$
31: **end while**
32: **return** \mathcal{S}

B Reproducibility

All experiments are conducted on a single NVIDIA
A100 GPU. Key fine-tuning settings are: epochs =
15, learning rate = 2×10^{-5} , batch size = 8, max
sequence length = 512, weight decay = 0.01, and
early stopping with a patience of 3 epochs.

For uncertainty estimation, we use MC dropout
with the number of stochastic forward passes set to
 $K = 10$. In RADS, we train a dueling DQN sam-
pler for 300 episodes with ϵ -greedy exploration,
decaying ϵ from 1.0 to 0.05 with a multiplicative
factor of 0.995. We use an experience replay buffer
of size 10000 and start network updates once at
least one minibatch is available (batch size = 64).
Both the online and target networks are optimized
with Adam (learning rate = 10^{-4}) and discount

Transfer Learning to CHIFIR Strategy		Performance on CHIFIR				Performance on PIFIR			
Datasets		Acc	F1	P	R	Acc	F1	P	R
Baseline	CHIFIR	0.942	0.824	0.778	0.875	–	–	–	–
Zero-shot	PIFIR	0.154	0.267	0.154	1.000	0.714	0.812	0.788	0.839
Full-shot	PIFIR + CHIFIR	0.904	0.615	0.800	0.500	0.857	0.900	0.931	0.871

Table 6: Performance comparison of transfer learning from PIFIR to CHIFIR under zero-shot transfer and full-shot transfer.

Knowledge Transfer from PIFIR to CHIFIR												
Strategy	Performance on CHIFIR					Performance on PIFIR					Transfer Gap	
	Accuracy	F1-score	Precision	Recall	ROC-AUC	Accuracy	F1-score	Precision	Recall	ROC-AUC	$\Delta F1$	95% CI
Random	0.838	0.167	0.250	0.125	0.716	0.752	0.791	0.916	0.729	0.837	–	–
Uncertainty	0.788	0.267	0.286	0.250	0.656	0.905	0.938	0.909	0.968	0.880	0.671	[0.387, 0.967]
Diversity	0.154	0.267	0.154	1.000	0.591	0.738	0.849	0.738	1.000	0.883	0.584	[0.409, 0.752]
LM-DPP	0.154	0.267	0.154	1.000	0.489	0.738	0.849	0.738	1.000	0.815	0.583	[0.409, 0.752]
TAGCOS	0.769	0.143	0.167	0.125	0.545	0.929	0.954	0.912	1.000	0.827	0.811	[0.529, 0.986]
BatchBALD	0.846	0.000	0.000	0.000	0.486	0.714	0.793	0.852	0.742	0.742	0.793	[0.654, 0.897]
RADS	0.865	0.632	0.545	0.750	0.858	0.881	0.921	0.906	0.935	0.865	0.289	[0.075, 0.599]

Table 7: Transfer learning performance from PIFIR to CHIFIR with 8 samples selected in CHIFIR under different sample selection strategies. $\Delta F1 = F1(\text{PIFIR}) - F1(\text{CHIFIR})$. CI = Confidence Interval.

factor $\gamma = 0.95$, and the target network is synchronized every 10 episodes. $\rho = 0.9$. The reward is defined as the prior-aware utility minus a diversity penalty computed from the nearest-neighbor distance in the mean log-probability space $\bar{\ell}(x)$, with $\lambda = 0.01$.

Dataset Split Each dataset is split into training, validation, and test sets (around 70%, 10%, and 20%), preserving the original class balance. Table 8 shows the number of positive and negative samples in each split.

Split	CHIFIR			MIMIC-CXR			PIFIR		
	Total	P	N	Total	P	N	Total	P	N
Train	196	27	169	320	124	196	135	92	43
Dev	35	5	30	74	29	45	24	16	8
Test	52	8	44	99	38	61	42	31	11

Table 8: Class distribution for CHIFIR, MIMIC-CXR, and PIFIR across train, development, and test sets. P = the number of positive reports, N = the number of negative reports.

C Robustness under Imbalanced Sampling

We evaluate the robustness of our sampling strategy under imbalanced sampling. For transfer learning from CHIFIR to PIFIR, We randomly select 5 samples from PIFIR with different positive to negative ratios. Each setting is repeated five times to obtain stable results. Figure 9 shows the results. The best performance occurs when the positive

to negative ratio is 1.00:0.00, which matches the class ratio selected by our method. This happens because CHIFIR has many more negative cases. Prioritizing positive target samples helps counter this imbalance and narrow the target-domain class distribution gap during transfer.

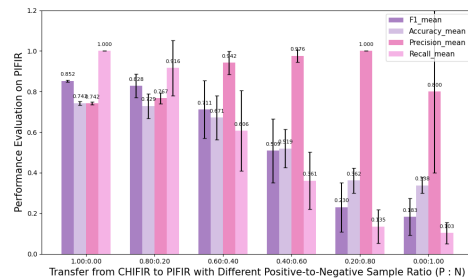


Figure 9: Class imbalance analysis of positive to negative sample ratios for CHIFIR to PIFIR transfer. Bars show mean values and black lines indicate variance.

D Transfer Performance from PIFIR to CHIFIR

Table 6 summarizes the baseline transfer performance. The zero-shot model trained on PIFIR exhibits severe performance degradation on CHIFIR, achieving very low accuracy and F1, indicating a substantial domain shift. In particular, the model tends to over-predict the positive class on CHIFIR (high recall but very low precision), which suggests that the decision boundary learned from PIFIR does not directly generalize to CHIFIR. When CHIFIR data are available for adaptation (full-shot), incor-

CHIFIR		PIFIR		MIMIC-CXR	
Word	TF-IDF	Word	TF-IDF	Word	TF-IDF
cells	18.405196	uptake	17.527355	chest	33.730506
fluid	12.311106	ct	16.972962	pneumonia	30.779783
bronchial	11.999560	fdg	14.968679	right	27.327142
description	11.998803	pet	11.530773	left	24.813217
biopsy	11.490593	right	9.951929	pleural	24.584867
tissue	10.625458	marrow	8.764460	effusion	22.224288
specimen	9.652244	activity	8.720572	pulmonary	21.700909
lung	9.026040	disease	8.699531	lung	21.272838
washings	9.022390	left	8.493931	comparison	21.251147
cell	8.798129	findings	8.308830	findings	20.950609

Table 9: Top 10 terms with the highest TF-IDF scores in CHIFIR, PIFIR and MIMIC-CXR.

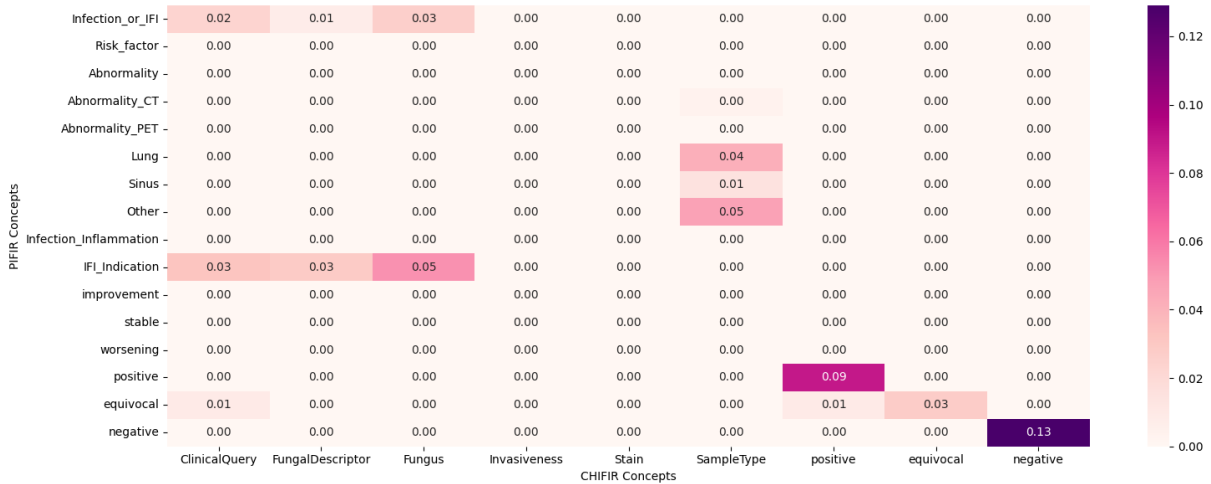


Figure 10: Jaccard similarity heatmap between CHIFIR and PIFIR concepts.

porating CHIFIR supervision mitigates this shift and improves target-domain performance, demonstrating that a small amount of target data is critical for reliable transfer.

We next evaluate whether our data selection strategy can maximize the benefit of limited target supervision. Table 7 reports transfer learning results when only 8 CHIFIR samples are selected for fine-tuning under different sampling strategies. Across all baselines, we observe that other strategies are insufficient to bridge the transfer gap: they either fail to improve CHIFIR F1 or produce unstable behavior. In contrast, RADS achieves the strongest transfer performance on CHIFIR, yielding the best overall target metrics (Accuracy/F1/ROC-AUC) while maintaining high performance on the source domain. Importantly, RADS also produces the smallest transfer gap ($\Delta F1$) among compared methods, indicating that the selected CHIFIR samples lead to more effective adaptation without sacrificing the knowledge learned from PIFIR. The confidence interval of $\Delta F1$ further suggests that RADS provides a more reliable and stable reduction of the trans-

fer discrepancy compared to alternative selection strategies.

E Learning Curves for MIMIC-CXR to CHIFIR Transfer under Varying Budgets

Figure 12 shows the performance from MIMIC-CXR to PIFIR across budgets under two baselines and Figure 13 (left) shows our methods’ performance. MIMIC-CXR is larger and the zero-shot baseline is already good. Therefore, the headroom for improvement is limited, and our method is similar to other baselines. Figure 13 (right) plots the domain gap $\Delta F1$ against budget with 95% confidence intervals. Across budgets, the point estimates stay close to zero and the confidence intervals largely overlap, indicating a small residual gap and no clear separation between budgets in this ultra-low-resource setting.



Figure 11: Concept-level KL divergence from CHIFIR to PIFIR.

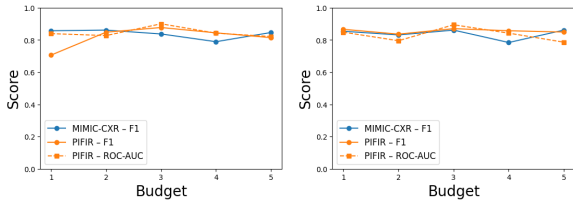


Figure 12: Transfer from MIMIC-CXR to PIFIR under baselines BatchBALD (left) and TAGCOS (right).

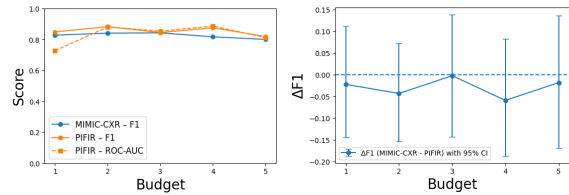


Figure 13: Transfer from MIMIC-CXR to PIFIR under our method RADS.

F Differences Analysis in CHIFIR, PIFIR, and MIMIC-CXR Datasets

CHIFIR contains 283 reports from 201 patients, with an average report length of 1,353 characters. PIFIR contains 201 reports from 156 patients, with an average report length of 1,809 characters. The MIMIC-CXR subset contains 493 reports from 290 patients, with a shorter average report length of 677 characters.

We first compute TF-IDF scores within each corpus and compare the top 10 highest-scoring terms between CHIFIR, PIFIR and MIMIC-CXR as shown in Table 9. CHIFIR is dominated by pathology and specimen-centric language (e.g., cells, fluid, bronchial, biopsy, tissue, specimen), reflecting cytology/histopathology reporting that em-

phasizes sample type and microscopic description rather than imaging observations. PIFIR is characterized by PET-CT and metabolic-imaging terminology (e.g., uptake, FDG, PET, CT, activity), as well as systemic disease descriptors (e.g., marrow, disease), consistent with PET-driven assessment of metabolic activity and whole-body involvement. MIMIC-CXR is dominated by chest radiography vocabulary and common pulmonary findings (e.g., chest, pneumonia, pleural, effusion, pulmonary, lung), reflecting the focus of X-ray reports on thoracic anatomy and acute cardiopulmonary abnormalities. Overall, these TF-IDF profiles highlight substantial modality- and workflow-driven lexical shifts between these datasets, motivating domain-adaptive transfer methods that can operate under pronounced vocabulary mismatch. Figure 3 also shows the transfer gap.

For the CHIFIR and PIFIR datasets, expert annotators also provided span-level annotation of concepts relevant to disease detection. Concept annotations in CHIFIR and PIFIR Datasets are listed in Table 12 and Table 11. The CHIFIR dataset reports 1,155 concepts, and the PIFIR dataset has 3,194 concepts. The two corpora serve different clinical niches. CHIFIR comes from cytology and histopathology notes and therefore focuses on microbiology terms such as FungalDescriptor and Stain. PIFIR is built from PET-CT reports and centres on imaging findings and risk factors, for example Abnormality_CT and Risk_factor.

To quantify overlap, we compute the Jaccard Similarity between the concept vocabularies:

Dataset	Example report excerpt (de-identified)
CHIFIR	""R groin LN biopsy"". Please note specimen has been received fresh and fragments have been sent to flow cytometry at 1225 on XXXXXX by XX/XXX. Two tan wispy cores 4 and 5 mm in length. A1. (dl:oze) ""R groin LN biopsy"". A tan core, 4mm in length with multiple fragments up to 1mm. A1. (dl/kr)
PIFIR	PET/CT technique: Scanning was performed encompassing the base of skull to upper thighs on a PET/CT scanner (GE 690 with time-of-flight). A contemporaneous low dose non-contrast multislice CT scan was performed for anatomic correlation and attenuation correction. Uptake time=70 minutes. BSL=<7mmol/L.
MIMIC-CXR	AP upright and lateral views the chest were provided. Mild left basal atelectasis. Lungs are otherwise clear. No signs of pneumonia or edema. No large effusion or pneumothorax. Cardiomedial silhouette is normal. Bony structures are intact. No free air below the right hemidiaphragm.

Table 10: Representative (de-identified) report excerpts from each dataset.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \quad (14)$$

where A and B are the sets of surface forms in CHIFIR and PIFIR, respectively. Figure 10 plots the resulting heatmap. Although both datasets include the classification terms positive, equivocal and negative, their lexical realizations share little common ground, so the Jaccard scores remain low.

Concept	Count	Unique	Diversity
Infection_or_IFI	279	174	0.62
Risk_factor	429	179	0.42
Abnormality	46	24	0.52
Abnormality_CT	460	204	0.44
Abnormality_PET	470	224	0.48
Lung	372	36	0.10
Sinus	19	4	0.21
Other	189	92	0.49
Infection_Inflammation	354	103	0.29
IFI_Indication	37	21	0.57
improvement	115	51	0.44
stable	29	16	0.55
worsening	55	34	0.62
positive	124	33	0.27
equivocal	129	68	0.53
negative	87	23	0.26

Table 11: Summary statistics for the IFI-related concepts in the PIFIR dataset.

Concept	Count	Unique	Diversity
ClinicalQuery	68	43	0.63
FungalDescriptor	294	86	0.29
Fungus	106	19	0.18
Invasiveness	39	27	0.69
Stain	172	16	0.09
SampleType	198	64	0.32
positive	118	40	0.34
equivocal	8	6	0.75
negative	152	12	0.08

Table 12: Summary statistics for the IFI-related concepts in the CHIFIR dataset.

To quantify directional lexical divergence, we compute the KL divergence on the concept level.

For each concept, P and Q denote the distributions of surface forms in PIFIR and CHIFIR, respectively. Both are smoothed over the combined vocabulary $\mathcal{V} = \mathcal{V}_{\text{PIFIR}} \cup \mathcal{V}_{\text{CHIFIR}}$ with a small ε to avoid zeros.

$$KL(P \parallel Q) = \sum_{v \in \mathcal{V}} P(v) \log \frac{P(v)}{Q(v)}, \quad (15)$$

where

$$P(v) = \frac{\text{count}_{\text{PIFIR}}(v) + \varepsilon}{\sum_{u \in \mathcal{V}} \text{count}_{\text{PIFIR}}(u) + \varepsilon |\mathcal{V}|}. \quad (16)$$

We compute $KL(\text{CHIFIR} \parallel \text{PIFIR})$ and visualize the results as a heatmap (Figure 11). In KL divergence, larger values indicate that greater mismatch between the two datasets. Even for shared classification terms (positive, equivocal, negative), the divergence values remain large. This suggests that the two datasets differ systematically in how their concepts are expressed, not simply in whether specific words occur.

Representative Report Examples To qualitatively illustrate domain- and modality-specific language, we provide representative (de-identified) report excerpts from each corpus in Table 10.