# Towards exploring continual learning for toxicologic pathology in pharmaceutical drug discovery

Anonymous CVPR submission

Paper ID \*\*\*\*\*

## Abstract

001 Machine learning has shown recent promise in advancing 002 toxicologic pathology image analysis during non-clinical safety assessment stages of drug development. However, 003 004 there exist notable challenges towards implementing these models in real-world scenarios, where datasets are se-005 quentially acquired over extended periods across weeks or 006 007 months in an Investigational New Drug study. A particular issue that arises in this process is the need for contin-008 uous adaptation to new data classes relevant to new tissue 009 010 types being scanned, or new animal models being incorporated, and so on. Data retention is often hindered by privacy 011 concerns, legal constraints, and storage limitations. Fur-012 thermore, existing deep networks are prone to catastrophic 013 forgetting when trained on new tasks, resulting in a sub-014 stantial loss of previously acquired knowledge. Therefore, 015 016 there is an urgent need for algorithms that are resilient to forgetting and capable of generalizing to new data without 017 018 the necessity of retaining large volumes of past examples or datasets. To address these challenges, we introduce a novel 019 020 replay methodology that leverages generative models, aug-021 mented by a knowledge regularization approach utilizing 022 attention embeddings from prior tasks. Our method inte-023 grates attention-based regularization, which prioritizes the 024 relative spatial importance of features, with generative latent replay. This synergistic approach enables the model 025 to retain and reinforce critical information from previous 026 027 tasks while adapting to new data. We empirically demonstrate the superior continual learning performance of our 028 029 method in non-stationary data environments, as evidenced by its application to a representative toxicologic pathology 030 image analysis task. 031

#### **1. Introduction**

A key component of non-clinical safety evaluations in phar maceutical drug discovery pertains to analysis of tissue
 morphologies in animal models subsequent to administra-

tion of the candidate drug in an Investigational New Drug 036 (IND) validation process. Toxicologic pathology is a spe-037 cialized field that integrates the study of toxicology (the 038 science of poisons) and pathology (the study of disease) to 039 understand and evaluate the effects of various agents on liv-040 ing organisms. This discipline is crucial in the regulatory 041 safety assessment process for predicting human and animal 042 responses to drugs, chemicals, and therapeutic devices, in-043 cluding their potential to cause adverse physiological and 044 functional effects. Toxicologic pathologists analyze tissue-045 level alterations to distinguish between incidental and test 046 article-related findings, crucial for accurate safety assess-047 ments in drug discovery. They identify alterations at the 048 cellular, subcellular, and molecular levels, considering fac-049 tors such as aging, genetics, and nutrition. This process 050 integrates toxicology and pathology to evaluate drug safety 051 and efficacy, predict human toxicities, and identify safety 052 liabilities early in development. Approximately 70% of the 053 toxicity-related attrition occurs in the preclinical phase [5], 054 making it essential to establish a preclinical safety margin 055 for viable drug candidates from the perspective of regula-056 tory guidance and the downstream feasibility of the candi-057 date. 058

The advent of deep learning has catalyzed significant ad-059 vancements in toxicologic pathology research, enabling un-060 precedented analytical capabilities. However, the integra-061 tion of these technological advances into clinical practice is 062 fraught with substantial hurdles, including the continuous 063 accumulation of large datasets over time and formidable 064 constraints related to privacy, storage, and data quality, 065 which impede the preservation of historical data. Contem-066 porary transfer learning techniques, despite their potential, 067 are often plagued by the phenomenon of catastrophic for-068 getting [4]. This occurs when a model's performance on 069 initially learned tasks dramatically declines as it is exposed 070 to new tasks, thereby compromising its reliability in deploy-071 ment environments and hindering its ability to adapt con-072 tinually to evolving data streams. Moreover, the storage 073 of even meticulously curated subsets of past samples is of-074 ten untenable due to stringent privacy regulations, logistical 075

impediments, and ethical considerations [14]. These challenges underscore the need for innovative approaches that
enable machine learning systems to learn continuously from
streaming data while preserving and leveraging knowledge
acquired from previous tasks.

081 In this paper, we introduce a novel, storage-free approach to continual learning in the analysis of animal histol-082 ogy images sourced from prior drug safety studies, designed 083 to overcome the aforementioned challenges. Our method-084 085 ology synergistically integrates two core components. We 086 employ a generative model to facilitate incremental-time replay from past data distributions. This model learns a la-087 tent space representation of previously encountered data, 088 enabling the sampling of synthetic data points that are in-089 terspersed with new data during the learning of subsequent 090 091 tasks. This strategy obviates the need for explicit storage of past data, thereby circumventing privacy and stor-092 age constraints. Concurrently, we preserve a snapshot of 093 the most salient regions from prior classes by constructing 094 class-specific attention embeddings. These embeddings en-095 096 capsulate the most informative spatial features pertinent to the model's decision-making process for a given task at time 097 point t. During the learning of new tasks at time t+1, these 098 attention embeddings are employed for regularization, en-099 suring that the model retains and leverages the most critical 100 aspects of previously acquired knowledge. 101

102 This dual framework introduces a novel pathway for integrating the most relevant regions into the model's 103 decision-making process when learning representations for 104 105 tasks at a given time point. The representations thus formed encapsulate a snapshot of previously encountered classes, 106 107 which is presented to the model alongside sampled data 108 points from the latent space learned via the generative model. While knowledge distillation methods have been 109 exploited in prior continual learning systems, the augmen-110 111 tation of learned representations of prior tasks using the most contextually important subsets of inputs, or a latent-112 113 space replay for past data distributions, remains largely unexplored. Drawing inspiration from the notion that inter-114 pretability techniques encode the subsets of the image space 115 most influential in model inference, we posit that prior-116 itizing and amplifying the preservation of such informa-117 118 tive spatial features during incremental learning sessions on novel classes or data can significantly enhance the model's 119 120 ability to retain the most vital subsets of previously seen data. Our approach is underpinned by a robust theoret-121 122 ical framework and is meticulously evaluated on a range 123 of digital pathology datasets, encompassing both histology 124 images and toxicologic pathology data. Through extensive experimentation, we demonstrate that our method not 125 only mitigates catastrophic forgetting but also enhances the 126 model's adaptability, robustness, and generalization capa-127 128 bilities in dynamic real-world clinical settings. Furthermore, we show that our approach outperforms state-of-the-<br/>art continual learning techniques, thereby establishing its<br/>potential as a vital tool for toxicologic pathology and be-<br/>yond.129<br/>130131<br/>132

In the subsequent sections, we delve into the intricacies 133 of our proposed methodology, elucidate the experimental 134 setup, present a comprehensive analysis of the results, and 135 discuss the broader implications of our work. We con-136 clude by highlighting avenues for future research, empha-137 sizing the potential of our approach to revolutionize con-138 tinual learning in data-constrained environments. Our ap-139 proach is found to improve continual learning for health-140 care, ensuring that models can adapt to new data without 141 compromising the knowledge gained from previous tasks, 142 thus maintaining high performance and reliability. 143

# 2. Prior Work

While digital pathology in the clinical realm has seen a 145 surge of research seeking to integrate various machine 146 learning approaches towards enhanced understanding of 147 pathology imaging data, such approaches in preclinical 148 safety evaluations have been relative nascent. There has 149 been an emerging interest in evaluating deep learning for 150 regulatory toxicologic pathology applications [22], with 151 some recent studies documenting the automated identifica-152 tion of histopathologic lesions in whole slide images (WSIs) 153 from regulatory toxicity studies, utilizing pixel-based man-154 ual annotations. Kuklyte et al. [11] employed various con-155 volutional neural network (CNN) architectures to detect and 156 quantify histopathologic lesions in toxicity studies. Their 157 methodology involved training algorithms with pixel-based 158 annotations of a substantial number of histopathologic le-159 sions across five different organs. A comparable approach 160 was adopted to identify proliferative lesions in the positive 161 control group of Tg-RasH2 mouse carcinogenicity studies 162 [16]. Shimazaki et al. [20] developed a customized U-Net 163 architecture trained with pixel-based annotations of seven 164 distinct liver lesions. In contrast, Zehnder et al. utilized 165 an unsupervised approach [23], leveraging generative ad-166 versarial networks (GANs) and autoencoder architectures 167 trained on normal liver samples to minimize the anomaly 168 score for normal data. Consequently, exposing the model to 169 abnormal data results in reconstruction and discrimination 170 failures, leading to elevated anomaly scores. With the emer-171 gence of foundation models for digital pathology, exten-172 sions towards self-supervised representation learning per-173 tinent to regulatory toxicity evaluations have been proposed 174 for the discovery of morphomolecular signatures [7], and 175 also for more general applications towards building down-176 stream tasks for preclinical safety assessments [6]. 177

Efforts to address steep declines in model validation performance with non-stationary learning schedules has seen the emergence of methods: a) based on preserving weights 180

241

like EWC [10], or freezing parameter subsets and progres-181 182 sively expanding for new data[17], which causes computational and memory overheads; b) knowledge-based regu-183 larization using distillation [13] for providing a glimpse of 184 185 learnt representations from past tasks to enable the model finetuning during incremental adaptation to account for new 186 and old tasks' knowledge; c) rehearsal methods that use 187 generative models to recreate batches of previously seen 188 189 data points [8]. Consolidation methods have struggled to scale to complex real-world datasets in clinical imaging [8], 190 191 and replay approaches often require a storage of original inputs from past classes or tasks in a buffer for batch-mixing 192 during new task learning. Recently, buffer-free methods 193 have been explored in medical imaging [15, 21], primarily 194 using stored features or separate replay modules, sharing of 195 196 which causes regulatory complexities in hospital systems [21]. At the time of writing, we came across a recent work 197 proposing a latent replay regime similar to ours[12], but 198 uses non-public datasets and focuses mainly on domain in-199 200 cremental learning without considering the role of attention 201 in regularization of future learning by prioritising salient features of past classes. 202

203 Most explainability methods image analysis seek to identify subregions of images most influential to specific 204 model decisions. Usage of these ideas in continual learn-205 ing regimes was explored in [2], but was limited to using 206 image-level combinations for creating past subsets. We pro-207 208 pose to directly use information from explainability workflows to enhance learnt task representations towards enforc-209 ing stronger distillation-driven regularization. To accom-210 plish a buffer-free replay of samples pertinent to prior tasks, 211 212 we retain a latent representation of the classes from past 213 tasks to allow incremental time sampling using a generative model as an alternative to replay using retained origi-214 nal samples from past learning stages. To our knowledge, 215 216 the question of continual learning for toxicologic pathology evaluations in drug discovery has not been explored despite 217 the clear existence of data acquisition regimes that are tem-218 219 porally spaced in the non-clinical stages of pharmaceutical 220 development. This is crucial from an image analysis standpoint as the deployment of machine learning for the preclin-221 222 ical discovery workflows requires adaptation of algorithms 223 over the weekly or monthly arrivals of data, as well as from 224 different tissue types and animal models in the course of an 225 Investigational New Drug evaluation.

Our work proposes a continual learning methodology for 226 227 the development of robust and adaptable models for toxico-228 logic pathology image analysis, with implications for im-229 proving the accuracy of safety assessments and enhancing the understanding of tissue responses to pharmacological, 230 chemical, or environmental agents. By mitigating catas-231 trophic forgetting and enhancing the model's ability to gen-232 233 eralize to new data, our approach has the potential to faciliate more widespread adoption of deep learning based234image analysis for regulatory toxicologica pathology work-<br/>flows, which would help accelerate non-clinical safety eval-<br/>uation stages in preclinical development and enable expe-<br/>dited timelines for drug development in pharmaceutical set-<br/>tings.236236236237236238237239239

# 3. Methodology

**3.1.** Datasets

We utilized the toxicologic pathology imaging dataset cu-242 rated by Serna et al. [19] to validate our continual learning 243 protocol, as it effectively models the practical requirement 244 of multi-task tissue recognition over extended periods using 245 the multi-magnification whole slide image (WSI) dataset 246 contained therein. The dataset comprises nine preclinical 247 rat studies with tissue processing conducted at three distinct 248 contract research organization (CRO) laboratories. All stud-249 ies were performed on Wistar Han rats and Sprague-Dawley 250 rats. The original acquisition protocols, adhering to Swiss 251 regulations and approved by the Cantonal Ethical Commit-252 tee for Animal Research, ensured ethical compliance and 253 rigorous experimental standards. The WSIs in this study 254 encompass organs both with and without histopathologi-255 cal findings, and include scans from Hamamatsu (ndpi) and 256 Aperio (svs) scanners. The studies were intentionally se-257 lected to incorporate significant staining variations, thereby 258 facilitating the development and validation of algorithms ro-259 bust to real-world settings where diverse scanning, staining, 260 and storage protocols may be encountered during investiga-261 tional new drug studies. The associated metadata for the or-262 gans included in each WSI were derived from internal pro-263 tocols specifying the grouping of different organs per WSI. 264

The resulting dataset comprises N=320 whole slide im-265 ages, featuring the following organs: liver, kidney, thyroid 266 gland, parathyroid gland, urinary bladder, salivary gland, 267 mandibular lymph node, and others (negative class). The 268 liver and kidney were prioritized due to their critical rele-269 vance in preclinical safety assessments. Additionally, or-270 gans frequently embedded with them, such as submandibu-271 lar lymph nodes and salivary glands (embedded with liver), 272 and urinary bladder (embedded with kidneys), were in-273 cluded. The thyroid and parathyroid glands were also incor-274 porated to evaluate the detection of small, closely situated 275 organs. To enhance the positive selection of desired organs, 276 various confounding organs from nine organ sets were in-277 cluded as a negative class for model training (class "other"): 278 adrenal glands, aorta, and ureters; lung and heart; stom-279 ach and intestine; skeletal muscle, sciatic nerve, mammary 280 gland, and skin; prostate and seminal vesicles; testis and 281 epididymis; eye and harderian glands; bone with bone mar-282 row; and spinal cord. This comprehensive approach ensures 283 a robust and versatile dataset for advancing regulatory tox-284



Figure 1. An abridged, representative toxicologic pathology and image analysis workflow in a pharmaceutical drug development context, with digitisation and deep learning approaches built in for preclinical safety evaluations.

285 icologic pathology image analysis. It is noted that the con-286 founding tissue presence with the main tissues under study are not treated as separate independent classes in the contin-287 ual learning settings proposed, since the simulation of such 288 a setup requires us to have multiple tasks that are tempo-289 290 rally spaced, but the tasks themselves need to include well-291 defined classes being studied per task. This is a shortcoming of incremental task learning protocols for continual adapta-292 293 tion to new data streams in toxicologic pathology studies, and can be alleviated potentially through compute-efficient 294 295 adaptations of online learning or similar approaches, which 296 would be a non-trivial exercise given the significant memory and computational overheads involved in processing of 297 298 large whole slide images typically involved in such pathology studies and the requirement for several sequential anal-299 300 ysis steps to be conducted over a large number of animal models per study. 301

### **302 3.2. Problem Definition**

We explore a setting where a model is to be trained in an M-stage protocol, where every stage is a classification task with classes as  $X_t = \{X_{t,i}\}_{i=1}^{Q_t}, t \in [1, M]$ , with each X being a class representing its samples (i.e.,  $x_t \in X_t$ ), and  $Q_t$  being the number of classes at stage t. At stage t, the classifier from previous stage t - 1 is incrementally optimized over the classes at current stage. The objective of the learning is that, given a small threshold  $\epsilon$ , after the t stage optimization, the reductions in the inference accuracies over validation sets from all previous stages  $D = \{D_1, ..., D_{t-1}\}$ 

meet the following criteria:

$$\forall d \in D, d < \epsilon$$

We design this as an incremental learning experiment with 303 four classes in the initial training stage and four subsequent 304 classes in the incremental stage  $(M = 2, Q_1 = Q_2 = 4)$ . 305 This study is modelled as a sequential learning task over sets 306 of classes as above, with a proportion of classes being learnt 307 as base classes during an initial training stage. Next, the re-308 maining classes are learnt as incremental classes in a sub-309 sequent learning stage, leading to a multistage learning sys-310 tem over a temporal interval. The base classes are the liver 311 (LI), kidney (KD), Thyroid gland (TG) and urinary gland 312 (UG). The incrementally learnt classes include parathyroid 313 gland (PG), salivary gland (SL), mandibular lymph nodes 314 (MLN) and Other classes (BG). The former are used to op-315 timize for the initial task (task 1) and the latter are used 316 to adapt the model trained over base classes for the incre-317 mental task (task 2), thus simulating a continual learning 318 scenario. In the first stage, where task 1 is performed, a 319 Resnet-50 derived model is used to conduct a classification 320 task for the initial set of classes. The result of the first stage 321 using a cross-entropy loss is  $p = \operatorname{softmax}(z) \in \mathbb{R}^{Q_1}$ , where 322 z is the set of logits. The classification loss in the first stage 323 is defined as: 324

$$L_C(y,p) = -\sum_{i=1}^{Q_1} y_i \cdot \log(p_i)$$
 (1) 325



Figure 2. Our overall pipeline for computing attention embeddings and creating the GMM library for the latent replay operations during incremental learning over subsequent dataset arrivals corresponding to the new tissue classes.

where  $p_i$  is the predicted probabilities of classes in the initial task,  $y_i$  is the corresponding ground truth represented as a one-hot encoding.

# 329 3.3. Latent Replay

330 Many continual learning (CL) approaches employ a replay buffer to selectively retain samples from previous tasks, 331 which are then mixed with data from current classes in a 332 process known as rehearsal. This technique is designed 333 334 to mitigate catastrophic forgetting, where the model loses the ability to perform well on previously learned tasks as it 335 adapts to new ones. However, storing actual samples from 336 past tasks can be problematic due to memory constraints 337 338 and potential privacy concerns, particularly in sensitive do-339 mains such as medical imaging. To address these chal-340 lenges, we propose a novel latent replay method that performs rehearsal without the need to store any actual sam-341 ples from previous tasks. Instead, our approach involves 342 sampling from a learned latent space that represents past 343 344 classes. This is achieved by constructing Gaussian Mixture 345 Models (GMMs) for each class during each training session. The GMMs capture the statistical distribution of the data, 346 allowing us to generate synthetic samples that closely re-347 semble the original data without retaining any of the actual 348 349 samples. In the context of deep learning models, the initial 350 layers are typically responsible for extracting low-level fea-351 tures from the input data, such as edges, textures, and simple patterns. These features are generally applicable across 352 a wide range of tasks and datasets. Through pre-training 353 on initial datasets, the weights of these early layers stabi-354 355 lize and can be effectively repurposed for various applica-

tions, including complex tasks like medical image process-356 ing. For instance, Srivastava et al. [21] demonstrated the 357 efficacy of transfer learning in medical imaging, where pre-358 trained models were fine-tuned for specific medical tasks, 359 achieving state-of-the-art performance. Conversely, the lay-360 ers closer to the classification head of the model are tasked 361 with extracting high-level, discriminative features that are 362 specific to the classes and tasks at hand. These layers re-363 quire fine-tuning to optimize the model's accuracy for the 364 target task. To leverage this hierarchical feature extraction 365 process, we propose to extract fine-grained features from an 366 intermediate layer of the model, which we designate as the 367 replay layer. 368

Our method involves training generators using the ac-369 tivations from this replay layer. Instead of storing raw 370 histopathology samples in a buffer memory, we only store 371 the generators. This approach significantly reduces mem-372 ory requirements and mitigates potential privacy violations, 373 as the generators do not contain any actual patient data. 374 During the initial training phase for each class, we retain 375 the layer preceding the last batch normalization layer of a 376 ResNet50 backbone as the averaged representation of every 377 class per task stage. This averaged representation serves as 378 a compact and informative summary of the class-specific 379 features, which can be used to generate synthetic samples 380 for rehearsal during subsequent training sessions. By do-381 ing so, we ensure that the model maintains its performance 382 on previously learned tasks while adapting to new ones, all 383 without the need to store any actual samples from past tasks. 384 A GMM is constructed using these representations per class 385 for a task session, similar to the protocol proposed by [12]. 386

After each task session, we learn session-specific GMMs 387 388 per class using aggregated features obtained during training from the ResNet50 backbone. For the t<sup>th</sup> session and 389 ith class, the generator includes K multivariate Gaussians 390  $\mathbf{g}_1, \mathbf{g}_2, ..., \mathbf{g}_K$  in the mixture with probability density func-391 tion (PDF) as: 392

$$p(\mathbf{f}_t) = \sum_{k=1}^{K} w_k \cdot \mathcal{N}(\mathbf{f}_t | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
(2)

393

where  $\boldsymbol{\mu}_k = \frac{1}{N_t} \sum_{n=1}^{N_t} \mathbf{f}_t^n$  and  $\boldsymbol{\Sigma}_k = \frac{1}{N_t} \sum_{n=1}^{N_t} (\mathbf{f}_t^n - \boldsymbol{\mu}_k) (\mathbf{f}_t^n - \boldsymbol{\mu}_k)^T$  Here,  $N_t$  is the number of samples in ses-394 395 sion t for a class,  $\mathcal{N}(\mathbf{f}_t | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  is the PDF of the  $k^{\text{th}}$  com-396 397 ponent with mean  $\mu_k$  and covariance  $\Sigma_k$ . The value of  $K \in [1, K_{max}]$  is estimated using Bayesian Information Cri-398 399 terion approaches in cluster analysis [3]. For each K, pa-400 rameters of the mixture are computed by an EM algorithm [1]. Finally, the fitted model  $B_t^i$  enables generation of ran-401 dom samples for session i and class t. Over sequential task 402 403 sessions, a library of session-specific GMMs are accumu-404 lated to allow the creation of representations from any past classes and sessions or task stages for replay. 405

#### 406 **3.4.** Attention Guidance Regularization

The optimization of neural architectures has predominantly 407 408 relied on the empirical risk minimization framework facil-409 itated by gradient-based weight adaptation, conceptualizing the learning dynamics as analogous to the hill-climbing 410 problem. In this paradigm, the reduction of loss, defined 411 412 as the discrepancy between actual and predicted labels, is akin to descending a gradient slope [10]. Theoretically, this 413 414 process can be interpreted as the movement of gradient values within a conservative field, suggesting that the transition 415 from a loss value  $L_1$  to a lower minimum  $L_2$  is independent 416 417 of the specific parameter changes during the process and is 418 solely contingent on the initial and final parameter sets [4].

419 In practical terms, this implies that during transfer learn-420 ing across new distributions, there is no inherent regularization governing how parameter configurations may adapt to 421 optimize for the new task. Conversely, it is computation-422 ally infeasible to control shifts in the majority of the net-423 work's parameters. A potential compromise involves prior-494 425 itizing input features that exert a disproportionate influence on model outcomes, thereby indirectly enforcing regular-426 ization on the most critical model parameters. This insight 427 forms the basis of our approach, which incorporates atten-428 429 tion embeddings into the continual learning process.

While the computation of such embeddings and their 430 integration into sequential learning processes has been 431 sparsely explored in the literature, we investigated a rela-432 tively straightforward paradigm. In this approach, we assess 433 the suitability of attention embeddings as auxiliary branches 434 435 appended to the standard pre-softmax logits retained from previous learning instances.

While the generative latent replay stage allows sampling 437 from a latent space relying only on the mean and variance of 438 the learnt representations, the incorporation of a notion of 439 feature importance towards relaying the most salient parts 440 of a data distribution is crucial towards enabling past tasks' 441 regularization during incremental training on new classes. 442 This is achieved by 'attentive distillation', where we re-443 inforce past task embeddings for distillations with task-444 relevant attention vectors. We use class activation map-445 ping [18] for interpreting the model outputs in terms of a 446 heatmap of most relevant regions of the images used for 447 prediction. Such a step is carried out for the top-25% of 448 the instances in classes in the initial task, in the interest of 449 a trade-off between the expressiveness of the model repre-450 sentation and real-time memory constraints. For individual 451 images per class in the training data subset, we use the at-452 tention module to compute corresponding activation maps. 453

The class activation heatmaps so obtained are then vec-454 torized corresponding to the image representations of such 455 instances as projected in the penultimate fully-connected 456 layers of the model. The attention embeddings are then 457 cast into the dimensionality of the image representation log-458 its obtained in pre-softmax layers using average-pooling 459 operations, and a dot product operation is performed be-460 tween the attention vectors and relevant image embeddings. 461 This enables boosting of the most salient subsets of the 462 image space in the final representation. Finally, the class-463 specific embeddings that reflect both the image representa-464 tion and corresponding activations are obtained for usage in 465 the distillation-based regularization steps. Considering the 466 initial task logits prior to the superimposition of the atten-467 tion embeddings as  $w_{old}$ , the inclusion of the attention em-468 bedding  $a_{old}$  results in the overall old class logits as  $z_{old}$ , 469 which are weighted and used in the distillation objective for 470 incremental new tasks. In subsequent sessions, a distilla-471 tion term is used in the objective to allow inclusion of past 472 knowledge in the optimization (y') are final class probabili-473 ties for new task classes prior to softmax operation): 474

$$L_{distillation}(z_{old}, y') = -\sum_{i=1}^{N} \operatorname{softmax}\left(\frac{z_{old}}{T}\right) \cdot \log(\operatorname{softmax}\left(\frac{y'_{i}}{T}\right)) \quad (3)$$

Logits and predictions are scaled with a temperature hy-476 perparameter T in a softening process. It helps reduce the 477 disparity between the class label with the highest confidence 478 score in the probability vector with respect to other class la-479 bels and helps better reflect inter-class relationships at the 480 representation learning stage. Considering the overall logit 481 vector for old classes, after weighting as  $z_{old}$ , class-specific 482 logits are weighted to obtain a sum of class-weighted logits 483 as: 484  $K_1$ 

$$z_{old} = \sum_{i=1}^{K_1} u_i \cdot z_i \tag{4}$$
 485  
486

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

Table 1. Accuracy (%) for task 1/stage 1 classes, after task 1 is trained for, and after task 2 is incrementally added in stage 2. The difference in accuracies on validation sets of task 1 classes represents forgetting on task 1 classes due to task 2 addition. AM - Attention Modulated, wKD - weighted Knowledge Distillation, FT- simple Finetuning, '+' indicates combined usage

Stage	Stage 1 (for task 1, T1)					Stage 2 (for task 2, T2)					ΔAcc
	LI	KD	TG	UG	Avg(T1)	LI	KD	TG	UG	Avg(T2)	T2-T1
Our(AM+wKD)	82.45	80.77	78.36	79.44	80.26	78.65	76.84	74.22	76.91	76.66	3.61
Our(AM+KD)	82.45	80.77	78.36	79.44	80.26	73.62	72.46	70.81	70.21	71.77	8.49
Our (KD)	82.45	80.77	78.36	79.44	80.26	66.37	61.78	59.34	62.51	62.50	17.76
Our (AM+FT)	82.45	80.77	78.36	79.44	80.26	65.03	58.14	56.77	60.31	60.06	20.20
Ours (FT)	82.45	80.77	78.36	79.44	80.26	45.11	37.95	36.53	40.91	40.12	40.14
LwF.ewc [9]	82.45	80.77	78.36	79.44	80.26	61.75	58.37	53.82	57.26	57.80	22.46
Priv [21]	82.45	80.77	78.36	79.44	80.26	65.47	63.04	56.71	61.07	61.57	18.69
PCL [12]	82.45	80.77	78.36	79.44	80.26	65.29	63.78	56.32	60.02	61.35	18.91
MT [15]	82.45	80.77	78.36	79.44	80.26	65.98	63.17	58.43	62.33	62.48	17.78
DGR [21]	82.45	80.77	78.36	79.44	80.26	63.18	59.82	56.74	60.12	59.97	20.29

(5)

487 The logits from individual classes  $z_i, i \in [1, Q_1]$  are cal-488 culated by averaging pre-softmax probability values (after sigmoid activation) for examples from each of  $Q_1$  classes. 489 The weights  $(u_1, u_2, \ldots, u_{Q_1})$  are computed as inverse 490 of class-specific accuracy on validation sets of the initial 491 492 classes. The idea is to boost logits from classes which are 493 inherently difficult to learn for the model (lower the classspecific accuracy, higher the class weight). This reduces 494 the disparity among classes in their contribution towards the 495 overall sessional representation vector to be saved as an im-496 print of stage 1 learning. Overall, the net incremental task 497 498 training objective for learning beyond initial sessions then becomes ( $\gamma = 0.5$ ): 499

$$L = \gamma L_{crossent} + (1 - \gamma) L_{distillation}$$

#### **4. Experiments and Discussions**

502 The experiment is split into two sequential tasks, labeled task 1 and task 2. The initial task proceeds with a cross-503 entropy loss and task 2, the incremental task utilizes a joint 504 505 loss with a cross-entropy and a distillation term with a ResNet-50 feature extractor. The pre-softmax layer gener-506 ates probability scores by a sigmoid operation. An 80:20 507 split is used for train:test split on datasets. Input images 508 are resized to 224x224 and a batch size of 25 is used with 509 510 a learning rate of 0.0001 and Adam optimization. task 1 models are trained for 250 epochs on a (N,label) set for all 511 N frames. In task 2, models are trained for 250 epochs 512 on (N', label, logit) tuples, where the logit is the attention-513 514 modulated version of logits obtained in pre-softmax layers 515 after the initial task. We set T = 5.0 after grid search in 516  $T \in [1, 10]$ . Two 32 GB Nvidia V100 GPUs, 512MB RAM were used for training using the ResNet-50 base models 517 (~24.8 million parameters, implemented in Python 3.7.1 and 518 519 Tensorflow 2.0), with an average training time of 102s per 520 epoch observed in both tasks. In terms of the computational

landscape, the primary usage of floating point operations521stems from the model training procedures over the datasets522from each of the tasks involved. The computation of the at-<br/>tention maps, and the generation of the class-level sampled523embeddings from the GMMs in the library defined per task<br/>session, account for a minority of the computational budget.526

#### 4.1. Results

To design the incremental learning setup in the context of toxicologic pathology, Task 2 is initiated with four incremental classes following the completion of Task 1, where four initial classes were trained. The top 25% of samples per class are processed through the attention modules to generate attention maps. These maps are then vectorized, pooled, and transformed into embeddings with dimensions comparable to those obtained for corresponding images processed through the model. The dot product between image embeddings and their corresponding attention embeddings is aggregated over selected images per class to derive class-level embeddings. These embeddings are integrated into the net representation used for regularization in the distillation loss terms during the second stage of training for Task 2.

The inclusion of the most informative regions from the top quartile of correct predictions in prior tasks, along with sampling from the generators in the Gaussian Mixture Model (GMM) library constructed for the initial data, ensures robust task and class-level knowledge availability during incremental training. Notably, Task 1 classes exhibit consistent accuracies across compared baselines and proposed methods, as identical backbones are employed in a pure classification task for the initial stage.

In Table 1, evidence of catastrophic forgetting is observed, marked by a sharp decline in Task 1 validation accuracy following sequential optimization over Task 2 data  $(\Delta Acc)$ , when no intervention is applied to mitigate it. The implementation of our combined latent replay with attention-informed embeddings to preserve prior knowl-

558 edge effectively mitigates such steep performance declines. This is evident when comparing relative reductions in ac-559 560 curacy ( $\Delta Acc$ ) for methods that employ end-stage distillation alone. The dynamic interplay of attention-modulated 561 562 weighted distillation with intermediate concatenations of generative representations significantly reduces accuracy 563 declines, as measured by validation performance on Task 564 565 1 holdout sets before and after training on new task classes. 566 We adapted recent buffer-free methods in clinical datasets to compare our joint approach of generative replay and at-567 568 tentive distillation in our class-incremental setting. Our method outperforms strategies such as feature retention and 569 generative replay alone, indicating a non-trivial influence 570 of enhancing the visibility of the most salient image-level 571 572 features for past tasks' logit vectors used for regularization 573 in incremental training phases for the new classes' data. We evaluated knowledge regularization settings involving 574 575 class-weighted and non-weighted configurations, including versions that utilize attention embeddings, with the former 576 577 demonstrating performance gains across attention and re-578 play configurations.

579 Interestingly, measures to address past task forgetting also improve classification accuracies on new classes (Fig. 580 3), likely due to enhanced regularization causing forward 581 582 transfer effects which are potentially a result of more op-583 timal initializations of the parameter spaces during the initial optimization task. Overall, incorporating insights from 584 spatial feature importance over past inputs, along with gen-585 erative latent rehearsal, shows clear gains in the model's 586 learning regime, enforcing generalization when optimized 587 588 on new tasks' data in toxicologic pathology applications.





# 589 5. Conclusion

590 We present novel methods for continual learning for regu-1407 safety applications in drug discovery, using a public 592 toxicologic pathology dataset. Using attention embeddings 593 with a generative replay of samples in concordance with 594 previous data distributions is found to enhance the represen-595 tative power of the most salient input subsets without requir-596 ing storage of actual prior instances over time, and shown to

mitigate forgetting during cross-distillation based continual 597 adaptations on incremental tasks. This work represents an 598 early effort towards examining the question of tissue-level 599 image analysis in drug development contexts through the 600 lens of continual learning as a means towards performing 601 efficient machine learning based drug safety assessments in 602 the light of temporally-spaced arrivals of datasets during the 603 toxicologic pathology process. In future, we will adapt our 604 methods on systems reliant on transformer architectures and 605 similar approaches, study computational complexities of 606 scaling attention and generation modules on larger datasets 607 and any resultant scaling properties that might become evi-608 dent in the process, and integrate multimodal data relevant 609 to histopathological evaluations in clinical and preclinical 610 pathology workflows, including but not limited to, multi-611 omics data and associated metadata that may contain salient 612 information about diagnostic states under study. 613

643

644

645

646

647

648

649

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

### 614 References

- 615 [1] Frank Dellaert. The expectation maximization algorithm.
  616 *College of Computing, Georgia Institute of Technology*,
  617 2002. 6
- [2] P. Dhar, R. V. Singh, K. C. Peng, Z. Wu, and R. Chellappa.
  Learning without memorizing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
  pages 5138–5146, 2019. 3
- 622 [3] Chris Fraley and Adrian E Raftery. How many clusters?
  623 which clustering method? answers via model-based cluster
  624 analysis. *The computer journal*, 41(8):578–588, 1998. 6
- [4] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and
  Y. Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013. 1, 6
- [5] Cornelis ECA Hop. Compound attrition at the preclinical
   phase. Attrition in the Pharmaceutical Industry: Reasons,
   Implications, and Pathways Forward, pages 46–82, 2015. 1
- [6] Guillaume Jaume, Simone de Brot, Andrew H Song,
  Drew FK Williamson, Lukas Oldenburg, Andrew Zhang,
  Richard J Chen, Javier Asin, Sohvi Blatter, Martina Dettwiler, et al. Deep learning-based modeling for preclinical
  drug safety assessment. *bioRxiv*, 2024. 2
- [7] Guillaume Jaume, Thomas Peeters, Andrew H Song, Rowland Pettit, Drew FK Williamson, Lukas Oldenburg, Anurag
  Vaidya, Simone De Brot, Richard J Chen, Jean-Philippe Thiran, et al. Ai-driven discovery of morphomolecular signatures in toxicology. *bioRxiv*, 2024. 2
  - [8] R. Kemker and C. Kanan. Fearnet: Brain-inspired model for incremental learning. arXiv preprint arXiv:1711.10563, 2017. 3
  - [9] H.E. Kim, S. Kim, and J. Lee. Keep and learn: Continual learning by constraining the latent space for knowledge preservation in neural networks. In *International Conference* on Medical Image Computing and Computer-Assisted Intervention. Springer, 2018. 7
- [10] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, and A. Grabska-Barwinska. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, pages 3521–3526, 2017. 3, 6
- [11] Jogile Kuklyte, Jenny Fitzgerald, Sophie Nelissen, Haolin
  Wei, Aoife Whelan, Adam Power, Ajaz Ahmad, Martyna
  Miarka, Mark Gregson, Michael Maxwell, et al. Evaluation
  of the use of single-and multi-magnification convolutional
  neural networks for the determination and quantitation of lesions in nonclinical pathology studies. *Toxicologic Pathol-*ogy, 49(4):815–842, 2021. 2
- [12] P. Kumari, D. Reisenbüchler, L. Luttner, N. S. Schaadt, F.
  Feuerhake, and D. Merhof. Continual domain incremental learning for privacy-aware digital pathology. *arXiv preprint arXiv:2409.06455*, 2024. 3, 5, 7
- [13] Z. Li and D. Hoiem. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
  40(12):2935–2947, 2017. 3
- [14] Arijit Patra, Yifan Cai, Pierre Chatelain, Harshita Sharma,
  Lior Drukker, Aris T Papageorghiou, and J Alison No-

ble. Multimodal continual learning with sonographer eye-<br/>tracking in fetal ultrasound. In Simplifying Medical Ultra-<br/>sound: Second International Workshop, ASMUS 2021, Held672in Conjunction with MICCAI 2021, Strasbourg, France,<br/>September 27, 2021, Proceedings 2, pages 14–24. Springer,<br/>2021. 2676

- [15] H. Ravishankar, R. Venkataramani, S. Anamandra, P. Sudhakar, and P. Annangi. Feature transformers: privacy preserving lifelong learners for medical imaging. In *MIC-CAI 2019: 22nd International Conference*, Shenzhen, China, 2019. 3, 7
- [16] Daniel Rudmann, Jay Albretsen, Colin Doolan, Mark Gregson, Beth Dray, Aaron Sargeant, Donal O'Shea D, Jogile Kuklyte, Adam Power, and Jenny Fitzgerald. Using deep learning artificial intelligence algorithms to verify n-nitroson-methylurea and urethane positive control proliferative changes in tg-rash2 mouse carcinogenicity studies. *Toxicologic Pathology*, 49(4):938–949, 2021. 2
- [17] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. arXiv preprint arXiv:1606.04671, 2016. 3
- [18] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 6
- [19] Citlalli Gámez Serna, Fernando Romero-Palomo, Filippo Arcadu, Jürgen Funk, Vanessa Schumacher, and Andrew Janowczyk. Mmo-net (multi-magnification organ network): A use case for organ identification using multiple magnifications in preclinical pathology studies. *Journal of Pathology Informatics*, 13:100126, 2022. 3
- [20] Taishi Shimazaki, Ameya Deshpande, Anindya Hajra, Tijo Thomas, Kyotaka Muta, Naohito Yamada, Yuzo Yasui, and Toshiyuki Shoda. Deep learning-based image-analysis algorithm for classification and quantification of multiple histopathological lesions in rat liver. *Journal of Toxicologic Pathology*, 35(2):135–147, 2022. 2
- [21] S. Srivastava, M. Yaqub, K. Nandakumar, Z. Ge, and D. Mahapatra. Continual domain incremental learning for chest xray classification in low-resource clinical settings. In *MIC-CAI Workshop on Domain Adaptation and Representation Transfer*, pages 226–238, 2021. 3, 5, 7
- [22] Oliver C Turner, Famke Aeffner, Dinesh S Bangari, Wanda High, Brian Knight, Tom Forest, Brieuc Cossic, Lauren E Himmel, Daniel G Rudmann, Bhupinder Bawa, et al. Society of toxicologic pathology digital pathology and image analysis special interest group article\*: opinion on the application of artificial intelligence and machine learning to digital toxicologic pathology. *Toxicologic Pathology*, 48(2):277–294, 2020. 2
- [23] Philip Zehnder, Jeffrey Feng, Reina N Fuji, Ruth Sullivan, and Fangyao Hu. Multiscale generative model using regularized skip-connections and perceptual loss for anomaly detection in toxicologic histopathology. *Journal of Pathology Informatics*, 13:100102, 2022. 2