

Structured description of the MONKEY challenge

SUMMARY

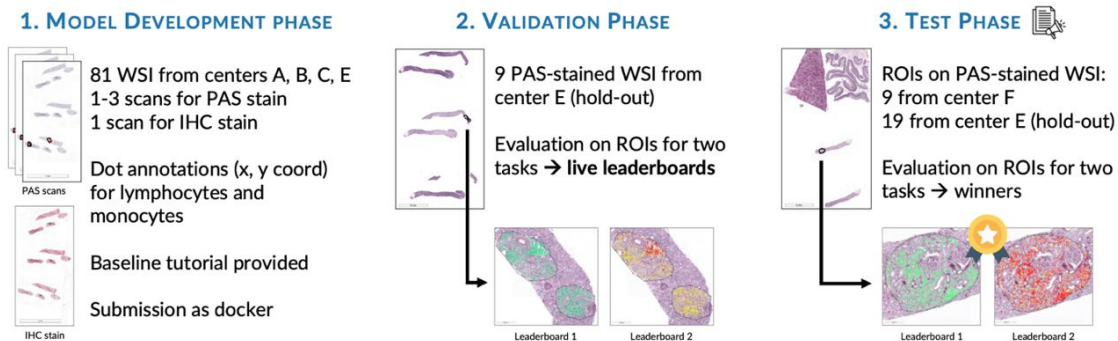


Figure 1: Overview of the three different phases of the MONKEY challenge. Phases one and two will run in parallel with a live leaderboard.

Item 1: Title

Machine learning for optimal detection of inflammatory cells in the kidney (MONKEY)

Item 2: Abstract

With a 2.5-fold increase in kidney transplants over the past 25 years, preserving the transplanted organ is crucial. To prevent transplant organ failure, early detection of risk factors and treatment is pertinent to patients' health. The Banff classification^{1,2} is the standard for histopathologic assessment of transplant kidney biopsies. It consists of 17 Banff Lesion Scores (BLS). Most BLS are graded semi-quantitatively as mild, moderate, and severe based on the number of inflammatory cells within the corresponding compartment. As the diagnosis and subsequent treatment decision depend on the result of the different BLS, it is of utter importance that the assessment of these individual BLS is objective and consistent.

Despite the Banff Classification System being very detailed with well-defined cut-offs for the different subcategories and a clear reference guide, it is not very robust due to its semiquantitative nature. It is difficult for pathologists to identify, label, and objectively and consistently grade each item, resulting in mild to moderate reproducibility scores (Kappa)³. Therefore, the development of automated BLS assessment holds great potential to reduce pathologists' workload and increase scoring consistency.

Since 8 of the BLS focus on the presence and extent of inflammatory cells in different kidney

¹ <https://banfffoundation.org/central-repository-for-banff-classification-resources-3/>

² Smith, Byron, et al. "A method to reduce variability in scoring antibody-mediated rejection in renal allografts: implications for clinical trials—a retrospective study." *Transplant International* 32.2 (2019): 173-183.

³ Schinstock, Carrie A., et al. "Banff survey on antibody-mediated rejection clinical practices in kidney transplantation: diagnostic misinterpretation has potential therapeutic implications." *American Journal of Transplantation* 19.1 (2019): 123-131.

compartments, we chose to address inflammatory cell detection with this challenge. Figure 1 shows an overview of the data, tasks and phases of the challenge. We provide dot annotations for mononuclear leukocytes (MNLs), which are the main type of inflammatory cells that are scored in the BLS. Previous research suggests that the subtypes of MNLs, monocytes and lymphocytes, play different roles in transplant rejection⁴, we also include the differentiation between the cells in the challenge. Differentiating these cells in the routine PAS staining can be tricky for pathologists, so it will be interesting to see how AI algorithms perform on this task. To provide a strong ground truth reference we used a specialized re-staining to be certain about the subtype of the inflammatory cells. To the best of our knowledge, this is the first publicly released dataset of its kind. Since we focus on the routine staining, the developed algorithms are also directly applicable for the diagnostic scoring.

Acknowledgements: This challenge is part of the DIAGGRAFT project⁵ at Radboudumc (Nijmegen, Netherlands) funded by the Dutch Kidney Foundation,

Item 3: Keywords

Kidney transplant biopsies, Cell detection, Inflammation detection

CHALLENGE ORGANIZATION

Item 4: Organizers

Core organization team:

- Linda Studer, Department of Pathology, Radboudumc, Nijmegen, The Netherlands
- Dominique van Midden, Department of Pathology, Radboudumc, Nijmegen, The Netherlands
- Luuk Hilbrans, Department of Nephrology, Radboudumc, Nijmegen, The Netherlands
- Jesper Kers, Department of Pathology, Amsterdam University Medical Centers, and Center for Analytical Sciences Amsterdam, Van 't Hoff Institute for Molecular Sciences, University of Amsterdam, Amsterdam, the Netherlands | Department of Pathology, Leiden University Medical Center, Leiden, the Netherlands
- Fazaal Ayatollahi, Department of Pathology, Radboudumc, Nijmegen, The Netherlands
- Jeroen van der Laak, Department of Pathology, Radboudumc, Nijmegen, The Netherlands

For additional contributors, see the website: <https://monkey.grand-challenge.org/organizers/>

Provide information on the **primary contact person**.

Linda Studer. Email address: linda.studer@radboudumc.nl

LinkedIn Profile: <https://www.linkedin.com/in/linda-studer/>

Item 5: Lifecycle type

This challenge will have two cycles.

1. **Challenge cycle:** This is the initial cycle. We will have a development and validation phase with a live leaderboard running for about 4.5 months, followed by a final test phase and an announcement of the winners.
2. **Open submission cycle:** after the announcement of the winners, the challenge will reopen for submissions and be supported for up to 5 years

⁴ Lamarthée, Baptiste, et al. "Transcriptional and spatial profiling of the kidney allograft unravels a central role for FcyRIII+ innate immune cells in rejection." *Nature communications* 14.1 (2023): 4359.

⁵ <https://www.computationalpathologygroup.eu/projects/diaggraft/>

Item 6: Challenge venue and platform

- a) Report the event (e.g., conference) that is associated with the challenge (if any).
We are applying to become a MIDL-associated challenge.
- b) Report the platform (e.g., grand-challenge.org) used to run the challenge:
The challenge is run on grand-challenge.org.
- c) Provide the URL for the challenge website (if any): monkey.grand-challenge.org

Item 7: Participation policies

- a) Define the allowed user interaction of the algorithms assessed (e.g. only (semi-) automatic methods allowed). **There is no user interaction.**
- b) Define the policy on the usage of training data. The data used to train algorithms may, for example, be restricted to the data provided by the challenge or to publicly available data including (open) pre-trained nets.
Additional data and pre-trained networks are allowed. The data sources must be reported, and either the data or the model weights must be publicly available under a permissive license.
- c) Define the participation policy for members of the organizers' institutes. For example, members of the organizers' institutes may participate in the challenge but are not eligible for awards.
Members of our institution are allowed to participate like regular contestants but are, of course, not permitted to access any of the test data on our local data shares. Members from the core organizing team are, however, excluded.
- d) **Define the award policy.** Provide details with respect to challenge prizes.
The total price money available is 3.250 EUR. It will be divided between the two leaderboards (monocyte and lymphocyte detection vs. mono-nuclear leukocytes (MNL), i.e., combined).
- e) **Define the policy for result announcement.**
Per the MIDL guidelines, we will hold a webinar at the end of the challenge, where the winners will be announced and can present their solutions. We will also publish a journal paper regarding our findings. Following the MIDL challenge guidelines, we will encourage all submitting groups to submit their method as a short paper to MIDL 2025.
- f) **Define the publication policy.** In particular, provide details on ...
... who of the participating teams/the participating teams' members qualifies as author
... whether the participating teams may publish their own results separately, and (if so)
... whether an embargo time is defined (so that challenge organizers can publish a challenge paper first).
Up to three members of each leaderboard's top three performing teams will be invited to participate in the challenge paper as consortium authors. Participants of the MONKEY challenge and non-participating researchers using the dataset can publish their own results at any time, separately. Challenge participants are encouraged to submit their solution as a short paper at MIDL 2025. Any such publications must cite this document (BIAS preregistration form for the MONKEY challenge), which will be updated on Zenodo (<https://zenodo.org/records/13382685>). Once a study protocol and/or a challenge paper has been published, they are requested to refer to those publication(s) instead.

Item 8: Submission method

- a) **Describe the method used for result submission.** Preferably, provide a link to the submission instructions.
Submissions will be made using Docker containers to grand-challenge.com. In addition

to the documentation by Grand Challenge, we provide instructions on the challenge page as well as a code tutorial on GitHub (<https://github.com/computationalpathologygroup/monkey-challenge>).

- b) Provide information on the possibility for participating teams to **evaluate their algorithms** before submitting final results.

There will be a live leaderboard phase, during which participants can see how their algorithm performs on a hold-out validation set of 9 cases. This should give them an indication on which algorithm to submit to the final test phase. They can also use cross-validation on the training dataset.

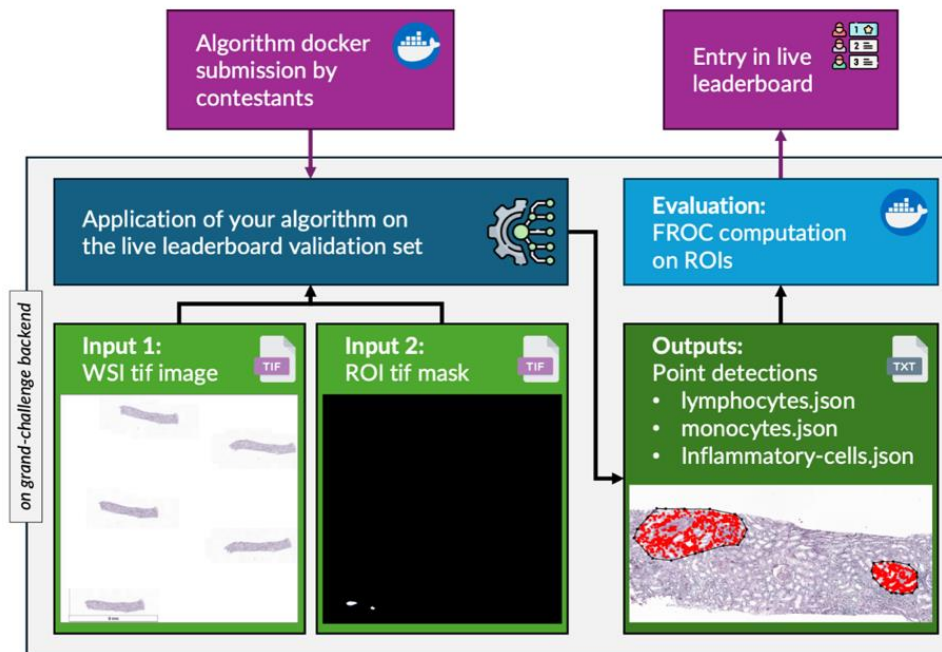


Figure 2: Overview of the submission workflow on grand-challenge.com for the live leaderboard phase. Contestants submit their detection algorithm as a docker to the challenge. The algorithm is then run on the backend on the hold-out validation set (i.e. the live leaderboard set). Then, the evaluation is run, and the metrics are reported back to the live leaderboard, which is displayed on the website.

Item 9: Challenge schedule

Provide a **timetable** for the challenge.

1. **Training data release** on the AWS Open Data Registry.
2. **Debugging phase** is opened: Contestants can submit containers which are evaluated on 2 slides from the training set. This allows us to release the error logs for easier debugging on the contestant's site.
3. **Kick-off Webinar**: We are holding an online webinar, which will be recorded and later released online, to introduce the challenge.
4. **Live leaderboard phase**: Submission to the validation phase is only possible when a team has made a successful submission to the debugging phase. The results of this phase will be immediately reported on a live leaderboard. We are using a hold-out set for this stage. Contestants can make 2 submission per week.
5. **Final test phase**: The contestants will have to choose which of their algorithm version to submit to the test phase, as there is only one submission possible. This phase will determine the winner of the challenge. We are using a hold-out set for this stage.

Timeline

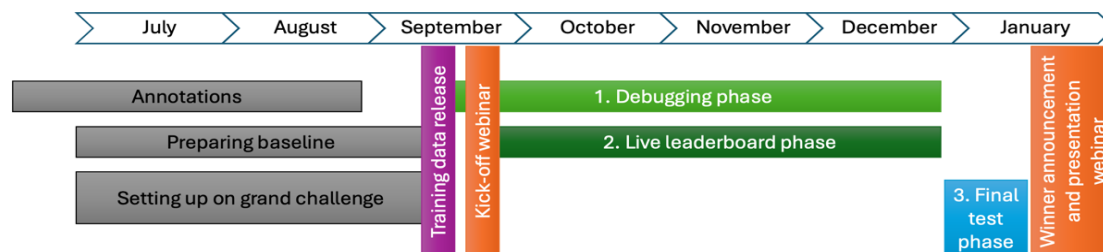


Figure 3: The debugging phase allows contestants to see the logs and test the functionality of their docker submission. The live leaderboard phase is the validation phase, where their algorithm is run on a hold-out set. Finally, there is a test phase with additional hold-out data to conclude the challenge winners.

Item 10: Ethics approval

Indicate whether **ethics approval** is necessary for the data. If yes, provide details on the ethics approval, preferably institutional review board, location, date and number of the ethics approval (if applicable). Add the URL or a reference to the document of the ethics approval (if available). Approval No. 2022-13686 from Prof. Dr. P.N.R. Dekhuijzen, Chair of the Institutional Review Board of the Radboud University Medical Center, CMO Radboudumc (METCoost-en-CMO@radboudumc.nl), approved on 31. March 2022. There are also data transfer agreements TAs in place with each data site that confirm their respective ethical approvals.

Item 11: Data usage agreement

Clarify how the data can be used and distributed by the teams that participate in the challenge and by others during and after the challenge. This should include the explicit **listing of the license** applied. The data will be distributed under the CC BY-NC-SA (Attribution-NonCommercial-ShareAlike) license.

Item 12: Code availability

- Provide information on the **accessibility of the organizers' evaluation software** (e.g. code to produce rankings). Preferably, provide a link to the code and add information on the supported platforms. The evaluation code will be provided to the participants on GitHub (<https://github.com/computationalpathologygroup/monkey-challenge>).
- In an analogous manner, provide information on **the accessibility of the participating teams' code**. The participating team's code and model weights must be available on GitHub (or a similar platform) and must be open access.

Item 13: Conflicts of interest

Provide information related to conflicts of interest. In particular provide information related to **sponsoring/ funding** of the challenge. Also, state explicitly who had/will have **access to the test case labels** and when.

This challenge is funded by the Dutch Kidney Foundation (Grant Nr. 21OK+012). The award money is a legacy from former IBEX employee John Theunissen. Any members of the Computational Pathology Team at Radboudumc can access the test case labels. The test cases and annotations will not be released publicly.

MISSION OF THE CHALLENGE

Item 14: Field(s) of application

State the **main field(s) of application** that the participating algorithms target.

- Diagnosis
- Medical image analysis research

Item 15: Task category(ies) State the task category(ies).

- Classification
- Detection
- Localization

Item 16: Cohorts

We distinguish between the *target cohort* and the *challenge cohort*. For example, a challenge could be designed around the task of medical instrument tracking in robotic kidney surgery.

While the challenge could be based on *ex vivo* data obtained from a laparoscopic training environment with porcine organs (challenge cohort), the final biomedical application (i.e. robotic kidney surgery) would be targeted on real patients with certain characteristics defined by inclusion criteria such as restrictions regarding gender or age (target cohort).

All data originates from transplant kidney biopsies collected from routine diagnostics. Thus, the data origin is the same for the challenge cohort and the target cohort. The cases were collected at different institutions, but the protocol for tissue preparation and PAS staining protocol is comparable and considered as routine practice. We also have different scanner modalities available (see Section 22.b.)), aiming to counteract challenges arising from image variations between the challenge dataset and a target cohort. The major difference between a target cohort and the challenge cohort is the distribution of different Banff categories and morphologies. To account for this, we collected an equal number of cases for different morphologies (listed in 22.b.)).

Item 17: Imaging modality(ies)

Specify the **imaging technique(s)** applied in the challenge.

For each case, there is a PAS-stained and IHC (double staining for CD3/CD20 and PU.1) (re-) stained whole slide image (WSI), performed in the lab at Radboudumc.

Item 18: Context information

Provide additional **information given along with the images**. The information may correspond...

- a.) ... directly to the **image data** (e.g. tumor volume).
- b.) ... to the **patient** in general (e.g. gender, medical history).

We will provide a quality score for the IHC slides, the institution from which the biopsy was taken, and the final biopsy diagnosis to give participants an overview of the distribution of different morphologies. The categories are insufficient clues for rejection (normal), ABMR (anti-body mediated rejection), TCMR (T-cell mediated rejection), mixed (ABMR+TCMR), borderline, chronic damage (IFTA), and other (BK virus nephropathy, necrosis).

Item 19: Target entity(ies)

- a.) Describe the **data origin**, i.e. the region(s)/part(s) of subject(s)/object(s) from whom/which the image data would be acquired in the final biomedical application (e.g. brain shown in computed tomography (CT) data, abdomen shown in laparoscopic video data, operating room shown in video data, thorax shown in fluoroscopy video). If necessary, differentiate between target and challenge cohort.

All data originates from transplant kidney biopsies, which are used to assess a transplant organ's health and define the treatment strategy that ensures the longevity of the donor organ. In the lab, the biopsies undergo FFPE (formalin- fixation and paraffin embedding) and are cut into thin tissue sections which are fixed onto glass slides. The glass slides are then stained with PAS (and IHC) and digitized using a WSI scanner.

Thus, both the challenge and target cohort consist of PAS-stained WSI.

- b.) Describe the **algorithm target**, i.e. the structure(s)/ subject(s)/ object(s)/ component(s) that the participating algorithms have been designed to focus on (e.g. tumor in the brain, tip of a medical instrument, nurse in an operating theater, catheter in a fluoroscopy scan). If necessary, differentiate between target and challenge cohort.

The challenge is focused on the development of algorithms for detecting and differentiating inflammatory cells in PAS-stained WSI. The inflammatory cells in question, also called mononuclear leukocytes (MNL), are monocytes and lymphocytes. The challenge is split into two leaderboards, 1.) MNL detection and 2.) detection of monocytes and lymphocytes (i.e., detection of MNL and differentiation between the two sub-classes). The number and distribution of these inflammatory cells needs to be scored by pathologists in routine kidney transplant diagnostics, thus the algorithm goal is directly transferrable from the challenge cohort to the target cohort.

Item 20: Assessment aim(s)

Identify the **property(ies) of the algorithms to be optimized** to perform well in the challenge. If multiple properties are assessed, prioritize them (if appropriate). The properties should then be reflected in the metrics applied (see "Metrics"), and the priorities should be reflected in the ranking when combining multiple metrics that assess different properties.

- *Example 1:* Find liver segmentation algorithm for CT images that processes CT images of a certain size in less than a minute on a certain hardware with an error that reflects inter-rater variability of experts.
- *Example 2:* Find lung tumor detection algorithm with high sensitivity and specificity for mammography images.

Corresponding metrics are listed below.

Our primary goal is to optimize the performance of inflammatory cell detection, followed by classifying the detected inflammatory cells into two subclasses. We will measure this using the Free Response Operating Characteristic (FROC) analysis (see "Metrics").

CHALLENGE DATA SETS

Item 21: Data source(s)

- a.) Specify the **device(s) used to acquire the challenge data**. This includes details on the device(s) used to acquire the imaging data (e.g. manufacturer) as well as information on additional devices used for performance assessment (e.g. tracking system used in a surgical setting).

See b.)

b.) Describe relevant details on the **imaging process/data acquisition** for each acquisition device (e.g. image acquisition protocol(s)).

All WSI are scanned with two different scanner settings (“CPG profile” and “diagnostic profile”) using a P1000 WSI scanner (3DHistech, Hungary) at Radboudumc. Different scan profiles result in different color reproduction in the resulting images, for which the developed AI models should preferably be insensitive. For most cases, we also offer the original scan performed at the source institution (Vienna: 3D-Histec Panoramic 250, Bern: 3D-Histec P1000, Emory: Olympus Nanozoomer, Mayo: Aperio system, Utrecht: Hamamatsu XR). We convert all WSI to TIF files with spacing of 0.24 mm/pixel. The average WSI file size is 490 GB, the average dimensions are 16 x 32 mm. For all cases, at least one ROI is annotated (total of 231 over the whole cohort). The ROIs have an average area of 0.32 ± 0.22 mm². There are an average of 350 lymphocytes and 180 monocyte point annotations per WSI.

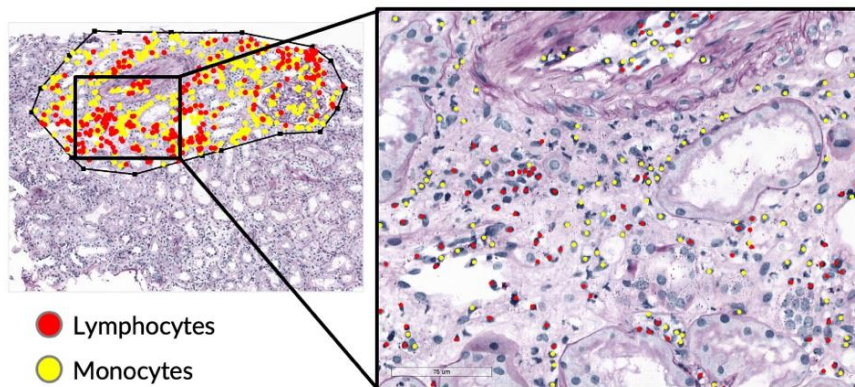


Figure 4: Example of annotated region of interest (ROI) with dot annotations for monocytes and lymphocytes (which combined are referred to as inflammatory cells).

c.) Specify the **center(s)/institute(s) in which the data was acquired and/or the data providing platform/source** (e.g. previous challenge). If this information is not provided (e.g. for anonymization reasons), specify why.

The cases and slides were collected six different pathology departments from 4 countries:

- A.) Radboudumc, Netherlands
- B.) UMC Utrecht, Netherlands
- C.) Medical University of Vienna, Austria
- D.) Mayo Clinic Minnesota, USA
- E.) IGMP, University of Bern, Switzerland
- F.) Emory University, USA.

d.) Describe **relevant characteristics** (e.g. level of expertise) of the subjects (e.g. surgeon)/objects (e.g. robot) involved in the data acquisition process (if any).

The inclusion criteria for selecting biopsies were set by the Banff digital pathology working group⁶. This criteria guideline was then provided the collaborating expert pathologists at the other data source institutions for case selection.

An experienced laboratory technician at Radboudumc developed the IHC retaining protocol, and renal pathologists evaluated staining quality and applicability. All retainings were performed at this lab. See also Item 23: Annotation characteristics.

⁶ The Banff Foundation for Allograft Pathology is an organization dedicated to advancing the understanding, research, and consensus development in transplant pathology, particularly focusing on the standardization of criteria for diagnosing allograft rejection in organ transplantation. See <https://banfffoundation.org/>.

Item 22: Training and test case characteristics

- a.) **State what is meant by one case in this challenge.** A case encompasses all data that is processed to produce one result that is compared to the corresponding reference result (i.e. the desired algorithm output).

Examples: Training and test cases both represent a CT image of a human brain.

Training cases have a weak annotation (tumor present or not and tumor volume (if any)) while the test cases are annotated with the tumor contour (if any). A case refers to all information that is available for one particular patient in a specific information (parameter 18). Both training and test cases are annotated with survival (binary) 5 years after (first) image was taken.

In the training set, one case consists of (i) 2-3 PAS-stained WSI scans: CPG profile, diagnostic profile, and original scan (if available), and (ii) double-stained IHC WSI scan. A case in the validation (live leaderboard) and final test set refers to one PAS-stained WSI scan (CPG profile).

The expected outputs are point predictions of the cells' location in millimeters.

- b.) State the **total number** of training, validation and test cases.

Training: $26+17+20+18 = 81$ (centers A, B, C, D)

Validation (live leaderboard): 9 (center E)

Test: 9 (center E) + 19 (center F) = 28

- c.) Explain **why a total number of cases** and **the specific proportion** of training, validation and test cases was chosen.

We split the cases by center to create a realistic scenario where a model is developed in a new institution. Four centers are used for the training set to ensure enough data was available for training. For these centers, all different stainings and scans will be made available for data augmentation purposes.

Center E is split between the validation and test set to have a reference in case certain algorithms perform well during validation and fail during test time. Center F is used as a hold-out set for the final testing. Both evaluation phases will be performed on the PAS slides scanned at Radboud with the CPG profile. See Table 1 for a more detailed overview.

Table 1: Overview of the case distribution between the different challenge stages. The PAS stained WSI are scanned in up to three different settings: Twice at Radboud with different scanning profiles (CPG and diagnostic), and once at the data source center. The slides in the training set will be made publicly available. The cases from Bern and Emory are used as hold-out sets for the validation (live leaderboard) phase and the final test phase (only the CPG profile).

	Radboud	Utrecht	Vienna	Mayo	Bern	Emory
Center label	A	B	C	D	E	F
Subset	train	train	train	train	val (9) test (9)	test
# Cases	26	17	20	18	18	19
# PAS – Radboud, CPG profile	26	17	20	18	18	19
# PAS – Radboud, diagnostic profile	14	17	20	18	18	19
# PAS original scan	n/a	17	n/a	18	18	19
# IHC slides	26	17	20	18	18	19

- d.) Mention **further important characteristics** of the training, validation and test cases (e.g. class distribution in classification tasks chosen according to real-world distribution vs. equal class distribution) and justify the choice.

Detection of inflammatory cells is strongly influenced by the presence of specific pathologies in the slides. Being relatively straightforward in biopsies with little

pathologies, it may be very hard (also for expert humans) to identify these cells in biopsies with severe scarring. We aimed to collect a similar number of cases and a similar distribution of morphologies from all centers. The split between validation and test for center E also ensures a similar distribution between both subsets.

The study protocol sent to all collaborators specified the following (aiming to collect 20 cases each):

- a. No-mild changes:
 - i. 2 no rejection or inconclusive
 - ii. 2 mild signs of rejection or mild IFTA (<25%)
- b. Moderate-severe changes
 - i. 2 moderate-severe glomerulitis
 - ii. 2 moderate-severe endovasculitis
 - iii. 2 moderate-severe tubulitis
 - iv. 2 moderate-severe peritubular capillaritis
 - v. 2 moderate IFTA (26-50%)
 - vi. 2 severe IFTA (>50%)
- c. Other alterations
 - i. 4 – tubulopathic changes, polyoma/BK, pyelonephritis, ischemic necrosis, etc.

Item 23: Annotation characteristics

- a.) Describe the **method for determining the reference annotation**, i.e. the desired algorithm output. Provide the information separately for the training, validation and test cases if necessary. Possible methods include manual image annotation, in silico ground truth generation and annotation by automatic methods.

We use IHC double staining for monocytes (applying monoclonal antibody PU.1, red) and lymphocytes (CD3/CD20, brown) to guide the annotation process. IHC involves selectively identifying cell-specific proteins (antigens). This is being visualized by using antibodies labeled with a chromogen/reagent. We used a CD3/CD20 double staining for lymphocytes, which results in a dark brown cytoplasmic staining. The PU.1 staining for monocytes results in a nuclear magenta staining. The IHC slides are re-stains of the PAS slide. An experienced laboratory technician developed the re-staining protocol, and renal pathologists evaluated staining quality and applicability.

- b.) If human annotation was involved, state the **number of annotators**.
6 annotators.

- c.) Provide the **instructions given to the annotators** (if any) prior to the annotation. This may include description of a training phase with the software. Provide the information separately for the training, validation and test cases if necessary. Preferably, provide a link to the **annotation protocol**.

The ROIs were all selected by DmV.

Student annotators were instructed on how to interpret the IHC staining to annotate the lymphocytes and monocytes by DvM. The annotation process was accelerated by generating automated detections of the lymphocytes and monocytes based on the IHC slide. The annotators (students + FM) then curated the automated annotations. False positive detections were deleted and missed monocytes and lymphocytes were added. The IHC restaining provides a very solid and easy to understand reference standard for annotations (see Figure 5), which allows non-pathologists to support the annotations process.

After this manual curation, DmV, who is a renal pathologist, reviewed all annotations. She also annotated cases with difficult morphologies or where the registration or

automated detection on the IHC failed. The ASAP software⁷ was used for all annotations. The same protocol was used for all cases.

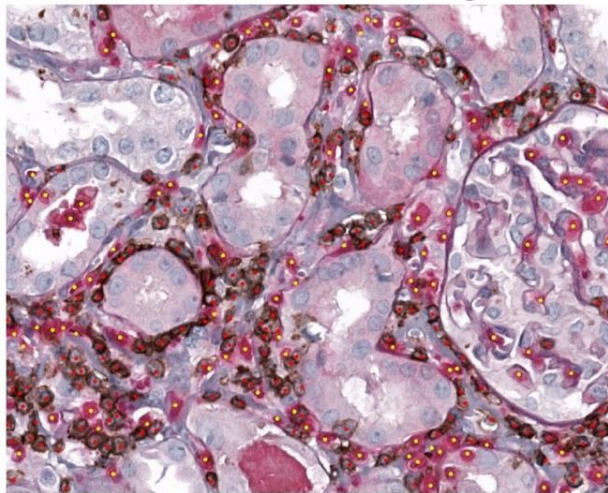
- d.) Provide **details on the subject(s)/algorithm(s) that annotated** the cases (e.g. information on **level of expertise** such as number of years of professional experience, medically trained or not). Provide the information separately for the training, validation and test cases if necessary.

We used HistoKat Fusion⁸ from Fraunhofer Mevis to align the original PAS to the re-stained slides. A previously developed model by Swiderska-Chadaj et al.⁹ for automated detection of lymphocytes was used to create automated annotations for both lymphocytes and monocytes (using color deconvolution). Five annotators, of which four are student annotators (VD, MdK, HQ, TdW) and one is a 5th-year resident pathologist (not specializing in renal pathology, FM),

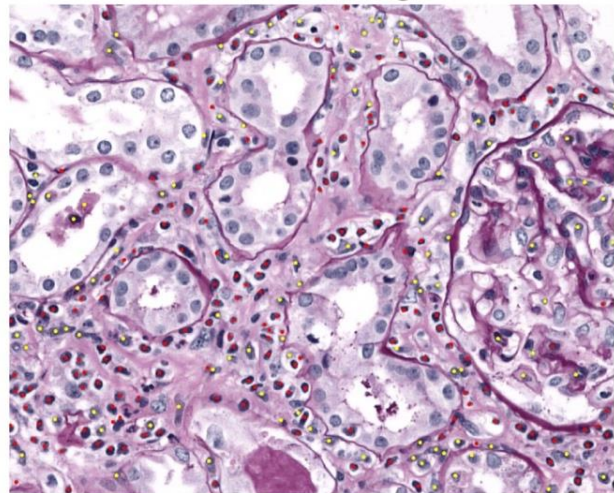
An expert renal pathologist (DvM) reviewed all annotations.

- e.) Describe the **method(s) used to merge multiple annotations** for one case (if any). Provide the information separately for the training, validation and test cases if necessary.
N/A

A. IHC double staining



B. PAS staining



● Monocytes ● Lymphocytes

Figure 5: Ground truth dot annotations overlaid on the (A) IHC-stained slide and the (B) PAS-stained slide. The lymphocytes and monocytes are well visible in brown (red dot annotations) and red (yellow dot annotations), respectively, on the IHC. Thus, the IHC staining provides an excellent reference for curating the automatically generated annotations and to transfer them to the PAS slide.

⁷ <https://computationalpathologygroup.github.io/ASAP/>

⁸ <https://histoapp.pages.fraunhofer.de>

⁹ Swiderska-Chadaj, Zaneta, et al. "Learning to detect lymphocytes in immunohistochemistry with deep learning." *Medical image analysis* 58 (2019): 101547.

Item 24: Data pre-processing method(s)

Describe the **method(s) used for pre-processing** the raw training data before it is provided to the participating teams. Provide the information separately for the training, validation and test cases if necessary.

Whole slide images (WSIs) are scanned using various staining techniques, profiles, and formats, necessitating their conversion to a standard format. All slides are first registered to the corresponding Periodic Acid-Schiff (PAS) staining diagnostic profile using HistokatFusion as the registration tool to achieve this. This ensures that all slides have the same coordinates, allowing annotations to be made and visualized on a common coordinate system. The output of the registration process is a file in the ".sqreg" format, which includes the paths to both the template and reference WSIs. All registered slides are subsequently converted to the TIFF format to standardize the format. The same protocol is used for all cases.

Item 25: Sources of error

- a.) Describe the most relevant possible error sources related to the image annotation. If possible, estimate the magnitude (range) of these errors, using inter-and intra-annotator variability, for example. Provide the information separately for the training, validation and test cases, if necessary.

In contrast to most other computational pathology challenges, MONKEY's reference standard is primarily based on highly objective immunohistochemical staining. Also, all annotations are visually checked and corrected if an experienced pathologist deems it necessary. Still, interpretation of IHC staining, varying IHC staining results, and manual operation of a computer mouse will introduce a small amount of annotation noise.

- b.) In an analogous manner, describe and quantify other relevant sources of error. Issues from scanning or WSI preparation (artifacts, bad staining, etc.) can negatively impact the image quality. However, since all ROIs are chosen manually, these issues are noticed and can be rectified, i.e., by rescanning or excluding the slide.

ASSESSMENT METHODS

Item 26: Metric(s)

- a.) Define the **metric(s) to assess a property of an algorithm**. These metrics should reflect the desired algorithm properties described in assessment aim(s) (parameter 20). State which metric(s) were used to compute the ranking(s) (if any).

We will apply the Free Response Operating Characteristic (FROC) analysis. In it, the true positive rate (TPR), a.k.a. sensitivity or recall) is plotted against the average number of false positives (FP) per mm² over all slides. We define a detected cell as true positive (TP) if it lies within a distance margin of a manually annotated cell. The margin is 10µm and 4µm for monocytes and lymphocytes, respectively. For the combined detection, the margin is 7.5µm. Based on this definition, we will compute the TP, FP, and false negatives (FN) and use them in the FROC analysis. From the FROC curve, we derive an "FROC score" by taking sensitivity at five pre-selected values of FP/mm²: [10, 20, 50, 100, 200, 300]. The score computation may be fine-tuned during the challenge to compare the best methods better.

- b.) **Justify why** the metric(s) was/were chosen, preferably with reference to the biomedical application.

FROC has been previously used in detection tasks in the TIGER and CAMELYON challenges.

Item 27: Ranking method(s)

- a.) Describe the **method used to compute a performance rank** for all submitted algorithms based on the generated metric results on the test cases. Typically, the text will describe how results obtained per case and metric are aggregated to arrive at a final score/ranking.
Leaderboard 1 will show the results for the MNL detection (overall inflammation). We will directly use the FROC value. Leaderboard 2 will show the lymphocyte and monocyte detection results, where the FROC will be computed per class and then averaged for the final ranking. The individual values will, however, be visible on the leaderboard.
- b.) Describe the method(s) used to manage submissions with missing results on test cases. Missing results will result in a lower performance, as there will be more false negatives. If a whole test slide fails, i.e. the result files are empty, the sensitivity will be zero as there will be no true positives resulting in a penalty to the overall sensitivity and thus the FROC.
- c.) **Justify why** the described ranking scheme(s) was/were used.
To address the class imbalance between the monocytes and lymphocytes.

Item 28: Statistical analyses

- a.) Provide **details for the statistical methods** used in the scope of the challenge analysis. This may include description of the **missing data handling**, details about the assessment of **variability of rankings**, description of any method used to assess **whether the data met the assumptions**, required for the particular statistical approach, or indication of any **software product** that was used for all data analysis methods.
We mitigate the risks and challenges outlined in Maier-Hein et al.¹⁰ as follows:
 - i. The case selection criteria were defined by the Banff digital pathology working group¹¹, which are a group of leading experts in the field of kidney transplant pathology
 - ii. The results are easily visually inspectable, as they are detections of individual cells for which we also have the IHC re-staining, which is a very reliable reference. We also have two pathologists in our core team that can give a visual assessment of the outputs. Additionally, we will use these algorithms in a reader study to see how pathologists and AI interact which will give us insight in how these detection algorithms can translate into routine diagnostics (see Item 29).
 - iii. We use the same metric as previous challenges with very similar tasks (TIGER challenge). We do not perform cross-metric aggregation, only case-based aggregation of the same metric (FROC).
 - iv. The evaluation code will be made publicly available on GitHub
 - v. All subsets are balanced regarding different morphologies (i.e. Banff rejection categories)
 - vi. We have a very robust ground truth thanks to being able to use IHC double staining as a reference.
- b.) **Justify why** the described statistical method(s) was/were used.
See above.

¹⁰ Maier-Hein, Lena, et al. "Why rankings of biomedical image analysis competitions should be interpreted with care." *Nature communications* 9.1 (2018): 5217.

¹¹ The Banff Foundation for Allograft Pathology is an organization dedicated to advancing the understanding, research, and consensus development in transplant pathology, particularly focusing on the standardization of criteria for diagnosing allograft rejection in organ transplantation. See <https://banfffoundation.org/>.

Item 29: Further analyses

After the final test stage of the challenge, we are planning to use the best algorithm(s) for future analysis:

- **Journal publication on challenge:** We will write a journal publication on the results the challenge. This will include the description of and reference of the best-performing algorithms and a more in-depth analysis and presentation of the results than provided by the leaderboard on Grand Challenge. Up to three co-authors per team will be invited as consortium co-authors (see Section 7. f).

All participants are also encouraged to submit at least a short paper to MIDL 2025 of their solution.

- **Susceptibility to image variance:** The original slides from both hold-out test sets (center E and F) have also been scanned with another scanner and another scanning profile (see Table 1, numbers indicated in italics). This allows us to assess the impact of scanner variability, especially since center F is using a very different scanner than most of the centers.
- **Reader study:** We are currently collecting an even larger cohort for a reader study with nine pathologists. The goal is for each of them to diagnose 100 cases. They will be requested to give scores for the six Banff lesion scores g, t, ptc, ci, ct, and i, as well as to categorize the biopsy into “T cell-mediated rejection”, “antibody-mediated rejection”, “mixed rejection”, “no specific allograft pathology”, and “other diseases of the allograft”. To compare the influence of computer-aided diagnostics, each pathologist will assess each case twice, once with and once without AI assistance. The AI generated results will be combination of the algorithms from this challenge, and a tissue segmentation algorithm that has previously been developed¹².

Thus, we can analyze the inter-observer agreement between pathologists, as well as between the pathologists and the AI.

Allograft failure and additional donor, recipient, and transplantation factors are also known for this cohort. This allows us to further study the prognostic value of pathologists’ diagnoses with and without AI assistance.

¹² Hermsen, Meyke, et al. "Deep learning–based histopathologic assessment of kidney tissue." *Journal of the American Society of Nephrology* 30.10 (2019): 1968-1979.