

VideoPipe 2022 Challenge: Real-World Video Understanding for Urban Pipe Inspection

Yi Liu^{1*}, Xuan Zhang^{1*}, Ying Li¹, Guixin Liang², Yabing Jiang², Lixia Qiu², Haiping Tang²,
Fei Xie², Wei Yao³, Yi Dai^{2†}, Yu Qiao^{1,4†}, Yali Wang^{1,5†}

¹ ShenZhen Key Lab of Computer Vision and Pattern Recognition, Shenzhen Institute of Advanced Technology,
Chinese Academy of Sciences, China

² Shenzhen Bwell Technology Co., Ltd, China

³ Shenzhen Longhua Drainage Co., Ltd, China

⁴ Shanghai AI Laboratory, Shanghai, China

⁵ SIAT Branch, Shenzhen Institute of Artificial Intelligence and Robotics for Society

Abstract—Video understanding is an important problem in computer vision. Currently, the well-studied task in this research is human action recognition, where the clips are manually trimmed from the long videos, and a single class of human action is assumed for each clip. However, we may face more complicated scenarios in the industrial applications. For example, in the real-world urban pipe system, anomaly defects are fine-grained, multi-labeled, domain-relevant. To recognize them correctly, we need to understand the detailed video content. For this reason, we propose to advance research areas of video understanding, with a shift from traditional action recognition to industrial anomaly analysis. In particular, we introduce two high-quality video benchmarks, namely QV-Pipe and CCTV-Pipe, for anomaly inspection in the real-world urban pipe systems. Based on these new datasets, we will host two competitions including (1) Video Defect Classification on QV-Pipe and (2) Temporal Defect Localization on CCTV-Pipe. In this report, we describe the details of these benchmarks, the problem definitions of competition tracks, the evaluation metric, and the result summary. We expect that, this competition would bring new opportunities and challenges for video understanding in smart city and beyond. The details of our VideoPipe challenge can be found in <https://videopipe.github.io>.

I. INTRODUCTION

In the last decades, sewer pipe system is one of the most crucial infrastructures in modern cities. In order to ensure its normal operation, we need to inspect pipe defects in an effective and efficient manner. Several technologies have been applied in the traditional pipe inspection procedure. [1] has conducted a thorough investigation and categorized them into visual methods, electromagnetic methods, acoustic methods, and ultrasound methods. In particular, Quick-View (QV) Inspection and Closed-Circuit Television (CCTV) Inspection are the most popular methods, as shown in Figure 1. The Quick-View (QV) Inspection is used for rapid anomaly assessment on sewer pipes, since the camera can only record videos on the pipe orifice. The CCTV inspection system involves a remote-controlled robot that travels along the sewer pipe with a camera for video recording [2]. Hence, it can get more

* Yi Liu (yi.liu1@siat.ac.cn) and Xuan Zhang (xuan.zhang1@siat.ac.cn) are equally-contributed first authors.

† Yi Dai (daiyi@bominwell.com), Yu Qiao (yu.qiao@siat.ac.cn) and Yali Wang (yl.wang@siat.ac.cn) are equally-contributed corresponding authors.

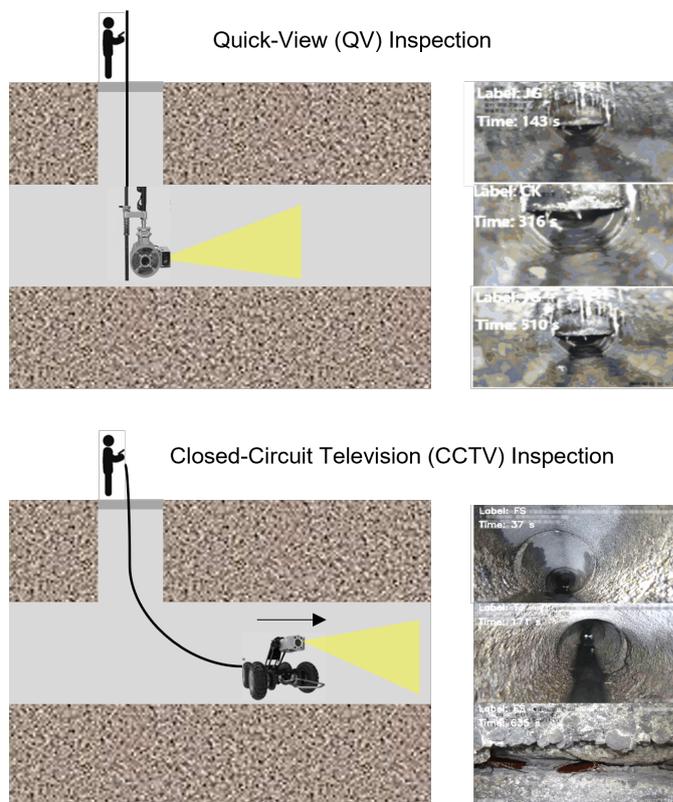


Fig. 1: Two widely-used pipe inspection methods.

detailed anomaly analysis for the whole pipe. Based on these QV and CCTV videos, the standardized protocols for manual inspection have been established and adopted in the recent years [3]. However, it is often labor-intensive to find anomaly from hundreds of hours of videos in the complex urban pipes.

To tackle this problem, it is essential to develop automatic inspection methods to discover sewer anomaly from large-scale pipe videos. Early works use hand-crafted visual features with traditional classifiers [4], [5], [6]. These approaches are often limited for inspecting defects in the complex scenarios.



Fig. 2: Anomaly Examples in Our QV-Pipe Dataset.

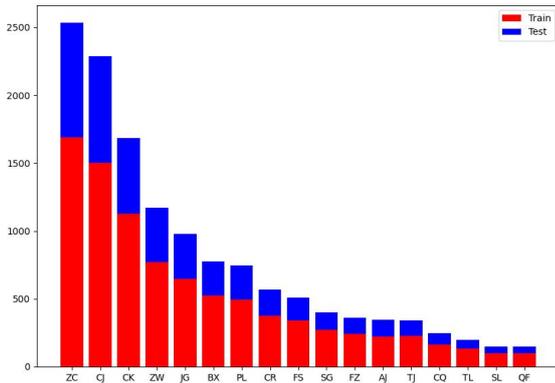


Fig. 3: Data Distribution of QV-Pipe Dataset.

With fast development of deep learning frameworks, a number of neural networks have been recently applied for defect classification [7], [8], [9], [10], [11], [4], [12], defect segmentation [13], [14], [15], and sewer defect detection [16], [17], [18], [19], [20], [21], [22]. However, these methods are mainly based on image-level detection. Hence, they are difficult to find multi-labeled and fine-grained defects, without learning spatial-temporal contexts in the video. Moreover, these data sets are not available for academic research, which blocks to develop advanced learning approaches for reliably inspecting sewer defects.

Based on the discussions above, we propose to establish real-world video understanding for urban pipe inspection. To achieve this goal, we carefully collect and annotate two new industrial video datasets, namely QV-Pipe and CCTV-Pipe, for our VideoPipe challenge. Specifically, QV-Pipe is used for video defect classification (Task 1) and CCTV-Pipe is used for temporal defect localization (Task 2). In this paper, we will introduce the details of data and tasks, and summarize the competition results. We expect that, this competition would bring new opportunities and challenges for video understanding in smart city and beyond.

II. DATASETS AND ANNOTATION

In this section, we present how these two datasets are collected and how they are constructed. Note that, these datasets are based on the real-world pipe networks. Hence, we have deleted the information of street, city and any other about privacy in our datasets.

A. QV-Pipe Dataset

The QV-Pipe dataset consists of 9.6k videos, which are collected from real-world urban pipes. The total duration of all videos exceeds 55 hours. Moreover, there are 1 normal class and 16 defect classes. Because the pipe situation is complex and multiple defects often appear at the same time, each video is annotated with multiple labels. The professional engineers are required to do this annotation procedure. To obtain accurate annotations of defect instances, these engineers are asked to

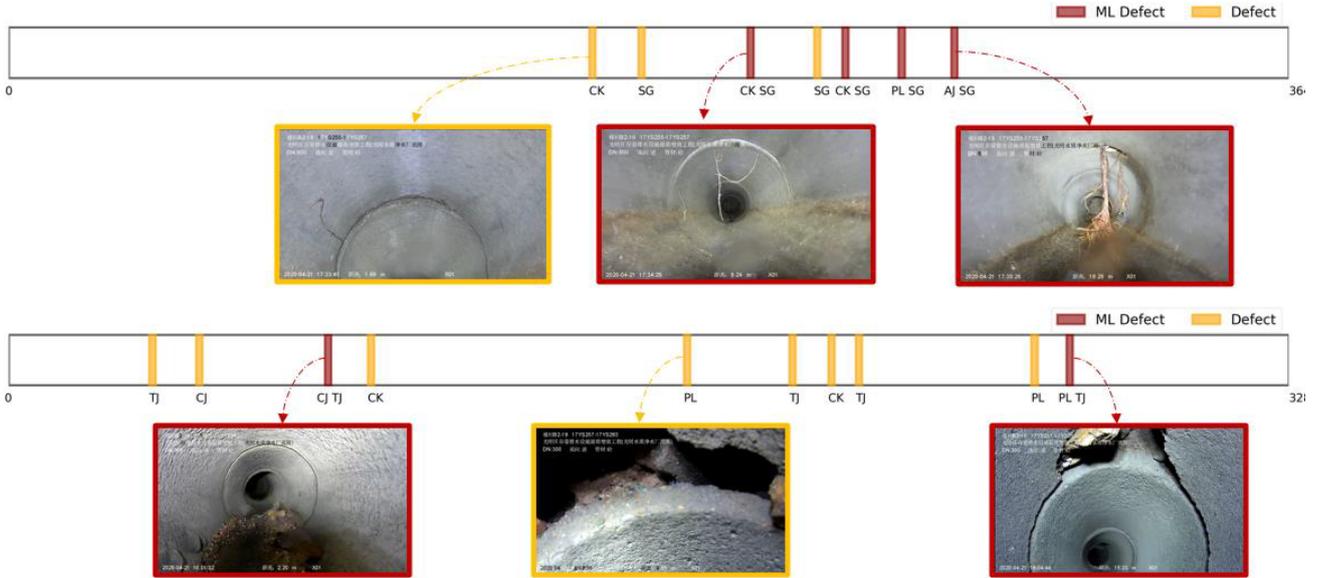


Fig. 4: Anomaly Examples of Our CCTV-Pipe Dataset. (ML: Multi-Labeled)

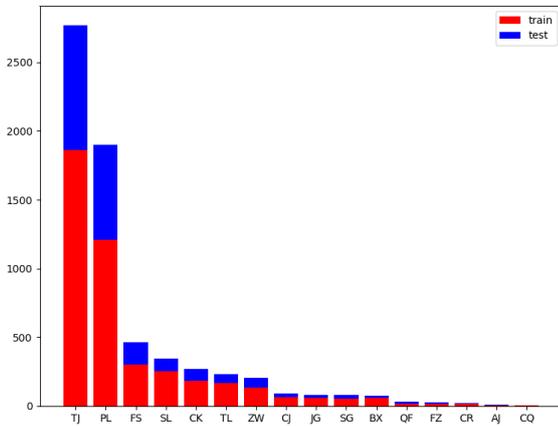


Fig. 5: Data Distribution of CCTV-Pipe.

check all the videos multiple rounds with cross validation. Examples of QV-Pipe are shown in Figure 2.

The QV-Pipe video duration ranges from 0.7 seconds to 385.2 seconds. Each video is annotated by 1 to 5 categories. On average, each video has the duration of 20.7 seconds and 1.4 labels. The 9.6k videos are divided into train set and test set according to the ratio of 2:1. As shown in Figure 3, the data exhibits the natural long-tailed distribution.

Moreover, we compare it with the existing benchmarks in video anomaly detection. As shown in Table I, our QV-Pipe dataset shows the following distinct characteristics. First, compared to the existing benchmarks, our QV-Pipe is large scale. Second, each video in our QV-Pipe contains multiple anomaly categories, and these categories are fine-grained. Finally, the previous datasets mainly works on human actions. Alternatively, the domain shift is large for urban pipe inspection.

Hence, our QV-Pipe brings new challenges and opportunities to understand video content for anomaly detection and beyond.

B. CCTV-Pipe Dataset

Our CCTV-Pipe dataset consists of 16 defect categories including structural and functional defects in the pipe. It contains 575 videos with 87 hours, which are collected from real-world urban pipe systems. Different from traditional temporal action localization, our goal in this realistic scenario is to find preferable temporal locations of defects from a untrimmed CCTV video, instead of exact temporal boundaries. Hence, professional engineers are asked to annotate a single frame for each defect. The annotation procedure has been checked multiple rounds with cross validation to guarantee label quality. We show some examples of CCTV-Pipe in Figure 4. We can see that, several defects appear at the same temporal location. Additionally, as demonstrated in Figure 5, the number of defects in each category ranges from 8 to 2,770. Such long-tailed distribution also raises new challenges for temporal defect localization.

Moreover, we compare it with the existing video benchmarks in temporal localization. As shown in Table II, our CCTV-Pipe dataset shows the following distinct characteristics. First, compared to the existing benchmarks, videos in our CCTV-Pipe can be very long in practice, e.g., average video duration is 545 s. It is quite challenging to find temporal locations of pipe defect from such long untrimmed videos. Second, instead of traditional segment annotation, we adopt single-frame annotation for realistic demand in urban pipe inspection. Moreover, multiple defects can densely appear at the same temporal location. These facts make our CCTV-Pipe as a challenging dataset for temporal localization. Finally, we compare it with the existing benchmarks in pipe defect inspection. As shown in Table III, our dataset is based on

Datasets	Multi-Labeled	Number of Classes	Number of Videos	Average Video Duration	Video Domain
UCSD Ped1 [23]	×	2	70	5 mins	Human Action
UCSD Ped2 [23]	×	2	28	5 mins	Human Action
Subway Entrance [24]	×	2	1	1.5 hours	Human Action
Subway Exit [24]	×	2	1	1.5 hours	Human Action
Avenue [25]	×	2	37	30 mins	Human Action
UMN [26]	×	2	5	5 mins	Human Action
RealWorld[27]	×	13	1,900	128 hours	Human Action
Our QV-Pipe	✓	17	9,601	55 hours	Pipe Defect

TABLE I: Video Anomaly Detection Benchmark Comparison.

Datasets	Multi-Labeled	Average Video Duration	Types of Video Annotation	Video Domain
THUMOS-14 [28]	×	261 s	Instance	Sports Action
ActivityNet [29]	×	117 s	Instance	Daily Action
HACS Segment [30]	×	149 s	Instance	Daily Action
Our CCTV-Pipe	✓	545 s	Single-frame	Pipe Defect

TABLE II: Temporal Localization Benchmark Comparison.

Datasets	Types of Data	Multi-Labeled	Number of Classes	Number of Images/Frames
Ye et al. [6]	Image-based	×	7	1 K
Myrans et al. [4]	Image-based	×	13	2 K
Chen et al. [7]	Image-based	×	5	18 K
Li et al. [10]	Image-based	×	7	18 K
Kumar et al. [9]	Image-based	×	3	12 K
Meijer et al. [11]	Image-based	✓	12	2,202 K
Xie et al. [12]	Image-based	×	7	42 K
Hassan et al. [8]	Image-based	×	6	24 K
Sewer-ML [31]	Image-based	✓	17	1,300 K
Our CCTV-Pipe	Video-based	✓	16	7,607 K

TABLE III: Urban Pipe Inspection Dataset Comparison.

videos, which is closer to urban pipe inspection in the real scenes. Moreover, our dataset is much larger than the existing ones, which opens new opportunities to develop powerful models for automatic defect inspection in urban pipe systems.

III. CHALLENGE TASKS AND EVALUATION

Video anomaly analysis is important for urban pipe system in the real world. Based on our QV-Pipe and CCTV-Pipe benchmarks, we introduce two challenging tasks in this competition, i.e., Video Defect Classification and Temporal Defect Localization, which aim at developing machine learning models to inspect pipe defects smartly and accurately.

A. Task Definition

Task 1: Video Defect Classification. Quick-View (QV) Inspection is one commonly-used technology. However, it is quite labor-intensive to find defects from a huge number of QV videos. To tackle this problem, we propose a video defect classification task, which is to predict the categories of pipe defects in a short QV video.

Task 2: Temporal Defect Localization. Closed-Circuit TeleVision (CCTV) is another popular method for pipe defect inspection. Different from short QV videos, CCTV videos are much longer and record more comprehensive content in the very distant pipe. The main task is to discover temporal locations of pipe defects in such untrimmed videos. Clearly, manual inspection is expensive, based on hundreds of hours of CCTV videos. To fill this gap, we introduce this temporal localization task, which is to find the temporal locations of

pipe detects and recognizing their corresponding categories in a long CCTV video.

B. Evaluation Metric

Each task has its own evaluation metric. Participants have been asked to upload the results according to the specified submission format. The submitted results have been evaluated according to different metrics.

Task1: Video Defect Classification. Since each video contains multiple categories, we use Average Precision (AP) to evaluate the recognition results on each defect category. Then we average AP over all the categories to obtain mAP.

Task2: Temporal Defect Localization. Referring to temporal action localization, we use Average Precision (AP) to evaluate the defect localization results on each defect category. Then we average AP over all the categories to obtain mAP. Due to our single-frame annotations, we compute temporal distance between the predicted defect and the ground truth to check if this prediction is a true positive. Finally, we use the average mAP as evaluation metric, which is the mean of mAP with all the temporal distances (from 5 second to 15 seconds, with 5 second interval).

C. Terms and Conditions

The videos of QV-Pipe and CCTV-Pipe are authorized by Shenzhen Bwell Robotics Co.,Ltd, which is one key member of our organization team. All these datasets are licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) license.

Task1	Team Name	mAP	Affiliation	Team Member
Top1	OverwhelmingFlt	72.18	Shanghai Paidao Intelligent Technology Co., Ltd., China	Jiawei Dong, Shuo Wang
Top2	fangxu622	61.99	Shenzhen University, China	Fang Xu
Top3	sthoduka	59.91	Hochschule Bonn-Rhein-Sieg, Germany	Santosh Thoduka

TABLE IV: Task 1 (Video Defect Classification): Top solutions that will be awarded.

Task2	Team Name	Avg. mAP	Affiliation	Team Member
Top1	OverwhelmingFlt	17.653	Shanghai Paidao Intelligent Technology Co., Ltd., China	Jiawei Dong, Shuo Wang
Top2	sthoduka	4.325	Hochschule Bonn-Rhein-Sieg, Germany	Santosh Thoduka

TABLE V: Task 2 (Temporal Defect Localization): Top solutions that will be awarded.

Hence, they must not be used for commercial purposes. Researchers can request access to these datasets by agreeing to terms and conditions in the following:

- The dataset is available for non-commercial research and educational purposes only.
- The users agree not to reproduce, duplicate, copy, sell, trade, resell or exploit for any commercial purposes, any portion of the images, and any portion of derived data.
- The users agree not to further copy, publish or distribute any portion of annotations of the dataset.
- We reserve the right to terminate your access to the datasets at any time.

For our challenge, the participants should also follow the terms and conditions below:

- The participants can use any pretrained models but cannot use any extra datasets for training in this competition.
- The participants should store the submission results for evaluation purposes.
- We request top-3 results to submit technical reports and/or code for us to verify submission validity.
- Any submitted reports and/or code will be used only for the sole purpose of evaluation.

IV. ORGANIZATION

Our challenge was run according to plan. First, we built the homepage (<https://videopipe.github.io>) for challenge introduction and released the datasets on February 28, 2022. Then, we held the challenges on the platform of CodaLab (Video Defect Classification: <https://codalab.lisn.upsaclay.fr/competitions/2232> and Temporal Defect Localization: <https://codalab.lisn.upsaclay.fr/competitions/2284>) from March 5, 2022 to May 5, 2022. The participants could only see the labels of training set. The results on test set are submitted through the platform and evaluated by the evaluation server.

V. SUBMISSION AND RESULTS

Overall, our challenge has attracted 89 participants for two tracks, where we received 51 and 38 participants respectively for Video Defect Classification and Temporal Defect Localization tracks. In Table IV and Table V, we list the top results of two tasks, which will be awarded in the competition. Both tasks clearly demonstrate that, it is necessary to hold this challenge, which encourages researchers and engineers to design more effective models for real-world video understanding. Note that, according to the competition rules, any teams

without technical report submissions will not be considered for an award. We will attach technical reports of these top results in our competition homepage, including method framework, dataset usage, training process, and inference process. In the following, we briefly summarize these solutions for each task.

Task1: Video Defect Classification.

1. *OverwhelmingFlt* (Shanghai Paidao Intelligent Technology Co., Ltd., China. Team Member: Jiawei Dong, Shuo Wang). A comprehensive solution of frame-based method, video-based method and super-image-based method. In the frame-based method, the participants use TRResNet [32] (ImageNet1K pretrained) in this task, and average prediction of all the frames as the final video prediction for inference. In the video-based method, the participants use Video Swin Transformer [33] (Swin-B backbone and Kinetics 400, Kinetics 600, or SomethingSomething V2 pretrained), and solve this task like traditional video classification. In the super-image-based method, the participants follow [34] to re-organize each video as a super-image, and use various image-based backbones (e.g., Convnext, TRResNet, NFNet, EfficientNet) for classification. Subsequently, the participants perform weight ensemble to summarize all the models for final prediction.

2. *fangxu622* (Shenzhen University, China. Team Member: Fang Xu). A solution with Video Swin Transformer [33] (Swin-B backbone). Furthermore, the participants use class-balanced focal loss to address class imbalance problem in the long-tailed dataset. The implementations are mainly based on MMAction2, with traditional settings and augmentations in video classification.

3. *sthoduka* (Hochschule Bonn-Rhein-Sieg, Germany. Team Member: Santosh Thoduka). A solution of ResNet-18 with a transformer encoder layer. In the training phase, 5 frames are sampled randomly from a train video. The focal loss is used to address class imbalanced problem in this task. In the test phase, 5 frames are sampled 5 times for a test video, and the final result is made on average of all the predictions.

Solution Trend: Overall, the solutions of participants are promising in the track of Video Defect Classification, based on the recent deep learning models with advanced training losses. From aspect of models, video backbones are often better than image backbones in this task. Especially, transformer-style operation can boost accuracy by learning complex long-range dependency among video frames. The super-image solution is also interesting in this task, which implicitly builds up frame relations via re-arranging a video as a super-image. From

aspect of losses, the focus loss is preferable to address class imbalanced problem in such as long-tailed data distribution.

Task2: Temporal Defect Localization.

1. *OverwhelmingFit* (Shanghai Paidao Intelligent Technology Co., Ltd., China. Team Member: Jiawei Dong, Shuo Wang). Basically, the participants attempt three temporal localization solutions, based on frame-level predictions, frame-level annotations and segment-level annotations. They choose the first solution due to its better performance. More specifically, it is a solution of frame-level prediction by ConvNeXt. First, a number of pseudo-label samples are generated to increase training set, since only single-frame annotations are available in this task. Second, they use ConvNeXt as image classifier for efficient training and testing. Finally, a number of post-processing approaches (e.g., moving average, peak finding, etc) are used to refine the final result from frame-level predictions.

2. *sthoduka* (Hochschule Bonn-Rhein-Sieg, Germany. Team Member: Santosh Thoduka). A solution of ResNet-18 with the addition of two linear layers is trained for the final prediction. Due to the limited training data, the participants enlarge the training set by frame sampling. Specifically, they select 15 frames in a window of $\pm 2.5s$ around each annotated frame as positive training samples, and randomly select extra normal frames as negative training samples. Moreover, they use a validation set to determine the best threshold for each defect class. Subsequently, they use these thresholds to distinguish if a frame is normal or contains defects in the testing phase.

Solution Trend: Compared to Video Defect Classification, Temporal Defect Localization is much more challenging. Since only single-frame annotations are given in this defect moment localization task, the participants tend to enlarge the training set, according to frame annotations. The post-processing step is also important to decide which frame may contain defects, due to complex scenarios in the real-world pipe system.

VI. CONCLUSION

This paper introduces the details of VideoPipe 2022 Challenge, which aims at building effective algorithms for urban pipe inspection in the real world. The challenge consists of two tasks including Video Defect Classification and Temporal Defect Localization, where we provide two large-scale video benchmarks (i.e., QV-Pipe and CCTV-Pipe) that are collected from real urban pipe systems. The results of this competition show that, machine learning algorithms can achieve promising results for pipe defect classification, but still have to be promoted for temporal defect localization. To sum up, this competition provides new opportunities and challenges for real-world video understanding. We expect that, in the coming future, these realistic video benchmarks would further show their impact on pattern recognition community and beyond.

VII. ACKNOWLEDGEMENT

This competition is partially supported by the National Natural Science Foundation of China (61876176,U1813218),

the Joint Lab of CAS-HK, the Shenzhen Research Program (RCJC20200714114557087), the Shanghai Committee of Science and Technology, China (Grant No. 21DZ1100100), Shenzhen Institute of Artificial Intelligence and Robotics for Society. It is also sponsored by Shenzhen Bwell Technology Co., Ltd, China.

REFERENCES

- [1] Z. Liu and Y. Kleiner, "State of the art review of inspection technologies for condition assessment of water pipes," *Measurement*, vol. 46, no. 1, pp. 1–15, 2013.
- [2] M. R. Halfawy and J. Hengmeechai, "Automated defect detection in sewer closed circuit television images using histograms of oriented gradients and support vector machine," *Automation in Construction*, vol. 38, pp. 1–13, 2014.
- [3] s. Rahman and D. Vanier, "An evaluation of condition assessment protocols for sewer management," 01 2004.
- [4] J. Myrans, R. Everson, and Z. Kapelan, "Automated detection of fault types in cctv sewer surveys," *Journal of Hydroinformatics*, vol. 21, no. 1, pp. 153–163, 2019.
- [5] Myrans, Joshua and Everson, Richard and Kapelan, Zoran, "Automated detection of faults in sewers using cctv image sequences," *Automation in Construction*, vol. 95, pp. 64–71, 2018.
- [6] X. Ye, J. Zuo, R. Li, Y. Wang, L. Gan, Z. Yu, and X. Hu, "Diagnosis of sewer pipe defects on image recognition of multi-features and support vector machine in a southern chinese city," *Frontiers of Environmental Science & Engineering*, vol. 13, no. 2, pp. 1–13, 2019.
- [7] K. Chen, H. Hu, C. Chen, L. Chen, and C. He, "An intelligent sewer defect detection method based on convolutional neural network," in *2018 IEEE International Conference on Information and Automation (ICIA)*. IEEE, 2018, pp. 1301–1306.
- [8] S. I. Hassan, L. M. Dang, I. Mehmood, S. Im, C. Choi, J. Kang, Y.-S. Park, and H. Moon, "Underground sewer pipe condition assessment based on convolutional neural networks," *Automation in Construction*, vol. 106, p. 102849, 2019.
- [9] S. S. Kumar, D. M. Abraham, M. R. Jahanshahi, T. Iseley, and J. Starr, "Automated defect classification in sewer closed circuit television inspections using deep convolutional neural networks," *Automation in Construction*, vol. 91, pp. 273–283, 2018.
- [10] D. Li, A. Cong, and S. Guo, "Sewer damage detection from imbalanced cctv inspection data using deep convolutional neural networks with hierarchical classification," *Automation in Construction*, vol. 101, pp. 199–208, 2019.
- [11] D. Meijer, L. Scholten, F. Clemens, and A. Knobbe, "A defect classification methodology for sewer image sets with convolutional neural networks," *Automation in Construction*, vol. 104, pp. 281–298, 2019.
- [12] Q. Xie, D. Li, J. Xu, Z. Yu, and J. Wang, "Automatic detection and classification of sewer defects via hierarchical deep learning," *IEEE Transactions on Automation Science and Engineering*, vol. 16, no. 4, pp. 1836–1847, 2019.
- [13] G. Pan, Y. Zheng, S. Guo, and Y. Lv, "Automatic sewer pipe defect semantic segmentation based on improved u-net," *Automation in Construction*, vol. 119, p. 103383, 2020.
- [14] C. Piciarelli, D. Avola, D. Pannone, and G. L. Foresti, "A vision-based system for internal pipeline inspection," *IEEE Transactions on Industrial Informatics*, vol. 15, no. 6, pp. 3289–3299, 2018.
- [15] M. Wang and J. C. Cheng, "A unified convolutional neural network integrated with conditional random field for pipe defect segmentation," *Computer-Aided Civil and Infrastructure Engineering*, vol. 35, no. 2, pp. 162–177, 2020.
- [16] X. Fang, W. Guo, Q. Li, J. Zhu, Z. Chen, J. Yu, B. Zhou, and H. Yang, "Sewer pipeline fault identification using anomaly detection algorithms on video sequences," *IEEE Access*, vol. 8, pp. 39 574–39 586, 2020.
- [17] S. Moradi and T. Zayed, "Real-time defect detection in sewer closed circuit television inspection videos," in *Pipelines 2017*, 2017, pp. 295–307.
- [18] S. Moradi, T. Zayed, F. Nasiri, and F. Golkhoo, "Automated anomaly detection and localization in sewer inspection videos using proportional data modeling and deep learning-based text recognition," *Journal of Infrastructure Systems*, vol. 26, no. 3, p. 04020018, 2020.

- [19] M. Wang, S. S. Kumar, and J. C. Cheng, "Automated sewer pipe defect tracking in cctv videos based on defect detection and metric learning," *Automation in Construction*, vol. 121, p. 103438, 2021.
- [20] J. C. Cheng and M. Wang, "Automated detection of sewer pipe defects in closed-circuit television images using deep learning techniques," *Automation in Construction*, vol. 95, pp. 155–171, 2018.
- [21] S. S. Kumar, M. Wang, D. M. Abraham, M. R. Jahanshahi, T. Iseley, and J. C. Cheng, "Deep learning-based automated detection of sewer defects in cctv videos," *Journal of Computing in Civil Engineering*, vol. 34, no. 1, p. 04019047, 2020.
- [22] X. Yin, Y. Chen, A. Bouferguene, H. Zaman, M. Al-Hussein, and L. Kurach, "A deep learning-based framework for an automated defect detection system for sewer pipes," *Automation in construction*, vol. 109, p. 102967, 2020.
- [23] W. Li, V. Mahadevan, and N. Vasconcelos, "Anomaly detection and localization in crowded scenes," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 1, pp. 18–32, 2013.
- [24] A. Adam, E. Rivlin, I. Shimshoni, and D. Reinitz, "Robust real-time unusual event detection using multiple fixed-location monitors," *IEEE transactions on pattern analysis and machine intelligence*, vol. 30, no. 3, pp. 555–560, 2008.
- [25] C. Lu, J. Shi, and J. Jia, "Abnormal event detection at 150 fps in matlab," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 2720–2727.
- [26] R. Mehran, A. Oyama, and M. Shah, "Abnormal crowd behavior detection using social force model," in *2009 IEEE conference on computer vision and pattern recognition*. IEEE, 2009, pp. 935–942.
- [27] W. Sultani, C. Chen, and M. Shah, "Real-world anomaly detection in surveillance videos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6479–6488.
- [28] H. Idrees, A. R. Zamir, Y.-G. Jiang, A. Gorban, I. Laptev, R. Sukthankar, and M. Shah, "The thumos challenge on action recognition for videos "in the wild";," *Computer Vision and Image Understanding*, vol. 155, pp. 1–23, 2017.
- [29] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles, "Activitynet: A large-scale video benchmark for human activity understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 961–970.
- [30] H. Zhao, A. Torralba, L. Torresani, and Z. Yan, "Hacs: Human action clips and segments dataset for recognition and temporal localization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8668–8678.
- [31] J. B. Haurum and T. B. Moeslund, "Sewer-ml: A multi-label sewer defect classification dataset and benchmark," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13 456–13 467.
- [32] T. Ridnik, H. Lawen, A. Noy, E. B. Baruch, G. Sharir, and I. Friedman, "Tresnet: High performance gpu-dedicated architecture," in *arXiv:2003.13630*, 2020.
- [33] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu, "Video swin transformer," in *arXiv:2106.13230*, 2021.
- [34] Q. Fan, C.-F. R. Chen, and R. Panda, "Can an image classifier suffice for action recognition?" in *ICLR*, 2022.