# AUTOMATED DEFECT DETECTION FOR MASONRY ARCH BRIDGES

*D. Brackenbury[1*], I. Brilakis[1] and M. DeJong[2]*

[1]*Department of Engineering, University of Cambridge. Cambridge, United Kingdom*
[2]*Department of Civil and Environmental Engineering, Berkeley University of California, Berkeley, USA*
*\* Corresponding author*

**ABSTRACT** The condition of masonry arch bridges is predominantly monitored with manual visual inspection. This process has been found to be subjective, relying on an inspection engineer's interpretation of the condition of the structure. This paper initially presents a workflow that has been developed that can be used by a future automated bridge monitoring system to determine underlying faults in a bridge and suggest appropriate remedial action based on a set of detectable symptoms. This workflow has been used to identify the main classes of defects that an automated visual detection system for masonry should be capable of detecting.

Subsequently, a convolutional neural network is used to classify these identified defect classes from images of masonry. As the mortar joints in the masonry are more distinctive than the defects being sought, their effect on the performance of an automated defect classifier is investigated. Compared to classifying all the regions of the masonry with a single classifier, it is found that where the mortar and brick regions have been classified separately, defect and defect free areas of the masonry have been predicted both with more confidence and with better accuracy.

## 1. Introduction

Currently the condition of masonry arch bridges is predominantly determined through manual visual inspections. This involves a detailed inspection conducted from touching distance of the structure, which is conducted at intervals of between two and twelve years, depending on the country. Detailed inspections are supplemented by superficial inspections which are conducted from ground level at intervals of between a half and one year (Helmerich et al., 2007). The manual visual inspection process is known to be subjective, as it is heavily dependent on the expertise and competence of the inspector. Laefer et al. (2010) assessed the reliability of the visual inspection process for detecting cracks in buildings by comparing the defects identified by two different inspectors. They found that there was a 14% difference in which cracks were detected by the inspectors, and that on average only 31% of cracks were identified. Furthermore, Phares et al. (2004) have performed a study which has demonstrated the variability of manual bridge inspections. Here 49 different bridge inspectors assigned a condition rating on a ten-point scale ranging from a failed condition to an excellent condition for seven different highway bridges. They have found that the different inspectors had on average assigned each element of the bridges between four and five different condition ratings out of the possible ten, showing that there is a large variation in the determined condition of the bridges between the different inspectors. From this study, they have also predicted that 78% of the average condition ratings assigned to bridges are incorrect at a 95% confidence interval. Automating defect detection and consequently the visual inspection process therefore has the potential to both increase the frequency and reduce the subjectivity of inspections.

The increasing capability and ease of geometric and photographic data acquisition presents an opportunity to create a digital visual model of bridges. However, this dataset alone has limited use unless it is augmented with information about defects on the structure and therefore the structural condition. Digital Imaging for Condition Asset Monitoring (DIFCAM) was a project in the UK looking to develop a capability for tunnel inspection to both capture and augment data with defect information. They created a road rail vehicle carrying an array of photographic sensors and a laser scanning sensor for data acquisition, as well as inertial and GPS sensors for position referencing. Digital Image Correlation was used to detect changes in the image and geometry data of the tunnel linings from one recording to the next. These changes represent defects that have developed in the structure between recordings (McCormick et al., 2014). This process is therefore unable to detect pre-existing defects, just changes in the structure from one recording to the next. Additionally, for the complex geometries of bridges, where lighting can't easily be controlled, it would be much harder to align the data taken between the two recordings to identify the changes.

Traditional approaches for directly detecting defects have relied on the assumption that defects will generally have a different intensity from the surrounding pixels. They have therefore looked to detect hand crafted features, such as sharp changes in pixel intensity or thresholding pixel intensity in order to detect defects.

More recent approaches have used machine learning to classify defects. These approaches can learn from diverse examples of defects, making them more robust. Samy et al. (2016) used a machine learning approach to detect defects in three dimensional images of masonry. They used a Support Vector Machine (SVM) to classify masonry images into different defect classes based on features extracted from them. The masonry images used were taken of a laboratory condition uniform brick wall with manually created defects. The image noise is therefore significantly less than in the case of a timeworn masonry arch bridge.

Deep learning approaches have the advantage over classical machine learning approaches in that they do not rely on handcrafted features for devising decision boundaries. Convolutional Neural Networks (CNNs) in particular have demonstrated state of the art performance for image classification tasks (Krizhevsky et al., 2012), and as a result they are the most popular network architecture for this purpose. Zhang et al. (2016) have compared the performance of a CNN to a SVM and a boosting method for detecting the presence of cracks in images of asphalt and have found the CNN to be superior with an $F_1$ score of 89.65%, 15% better than the other methods tested. Similarly Cha et al. (2017) have used a CNN to detect the presence of cracking in images of concrete with an accuracy of 98%. Chaiyasarn et al. (2018) have applied a CNN to images of masonry to detect cracking in image patches. They have achieved an accuracy of 74.9% but suggest that in some cases the system confuses the mortar joints with cracks.

Much of the focus of existing literature for defect detection has been on concrete and road surfaces. These experience many of the same defect classes as masonry, so there is potential for similar techniques to be applied to masonry as developed for concrete and road surfaces. However, Koch et al. (2015) reviewed different defect detection methodologies and concluded that the performance of defect detection algorithms with noisy data is questionable. Masonry images are inherently significantly noisier than concrete or road surface images due to the mortar joints between the individual masonry units. These mortar joints are often the most distinctive feature of masonry images, more so than the defects being sought. This is demonstrated by McRobbie (2009), who has attempted to apply a technique developed for concrete image surfaces to masonry image surfaces. He uses the Haar transform and image entropy to classify regions of images into those containing defects and those not. Whereas reasonable success was shown for concrete surfaces, with masonry surfaces the bricks and mortar have completely swamped any detected features.

This paper therefore investigates the effect of mortar joints on the performance of automated defect detection in masonry by comparing the detection accuracy where the mortar joints have been separated from the masonry images and tested for defects separately to that with no mortar joint separation. This comparison therefore determines the benefit of applying a two-stage methodology for detecting defects in masonry; first detecting and segmenting mortar joints, and then detecting defects. The classification is made using a state-of-the-art CNN classifier to detect defects in the presented images, therefore determining the applicability for similar techniques to those developed for concrete and asphalt road surfaces to be used on masonry. This is investigated for the detection of the different defect classes that are the most important to detect to determine the serviceability of masonry arch bridges.

## 2. Method

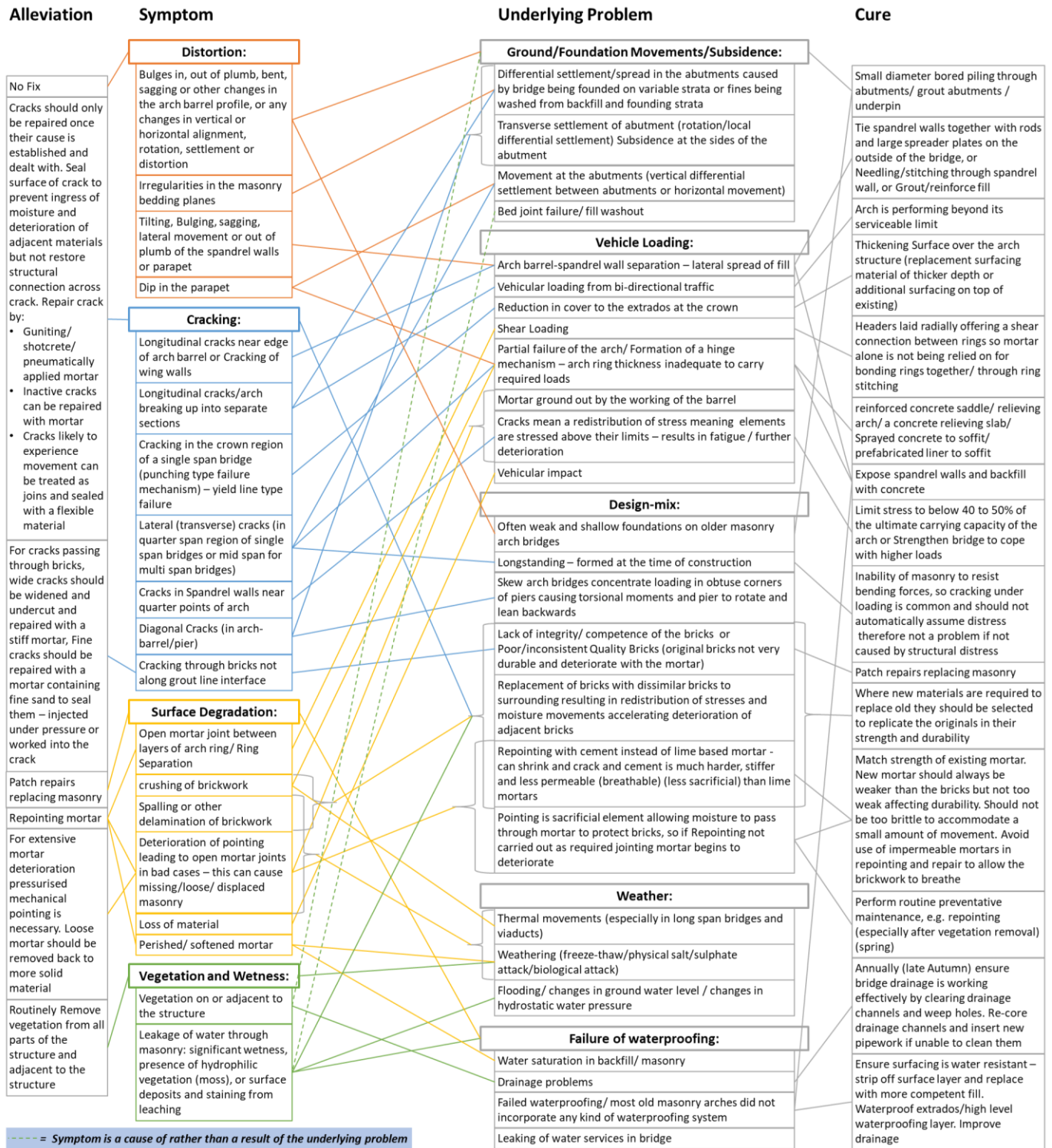### 2.1 Defect classes and problem identification

Information from the CIRIA documentation for assessing the condition of masonry arch bridges (McKibbins et al., 2006) and Network Rail standards for the examination of structures (Network Rail, 2017) has been linked to determine the severity and cause of the different defects that are visible on masonry arch bridges. The summary of this is shown in Figure 1, in which the different defects that are visible on the structure are described under symptoms, and these are linked to their root causes – the underlying problem that is causing the visible defect. The different repair strategies that are available both for correcting the identified underlying problems with the bridge and for repairing the visible defects on the bridge are then identified. In this way the visible defect on the bridge is linked to both its underlying problem and its solution. A future automated asset management tool for masonry arch bridges could therefore use the information in Figure 1 to determine the underlying problem in a bridge, and its required remedial treatment based on the detected defects.

The main classes of visible defects are; distortions in the shape of the bridge, irregularities in the mortar joints, cracking, spalling and other delamination of brickwork, missing or displaced masonry, mortar loss, vegetation, wetness, and surface deposits. Distortions in the shape of the bridge would be more accurately detected through examining the geometry by laser scanning than by visually examining the surface, and therefore this defect class has not been considered. The same is also true for missing masonry, though this visually resembles mortar loss or brickwork delamination, depending on the extent of missing masonry. As a result, this paper focuses on the visual detection of; cracking, spalling, mortar loss, and vegetation in images of masonry surfaces.

### 2.2 Dataset generation

A dataset of bridge images has been generated by closely photographing nine multi span masonry arch bridges near Cambridge. These bridges show widely differing masonry condition and appearance. This has generated approximately 24,500 images of masonry. From this dataset, 94 images have been chosen based on their depiction of defects. These images have had perspective distortion corrected, so that the masonry in the image is parallel to the image plane, generating image textures of the masonry surface. Since it is envisioned that defect detection will be performed on image textured three dimensional models of masonry arch bridges as part of an automated bridge inspection process, image textures of masonry surfaces are the expected input of a defect detection
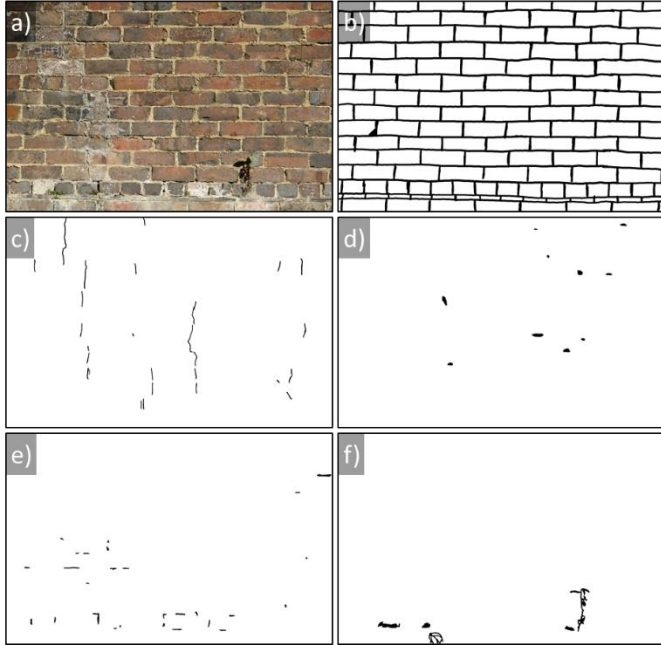
Figure 1 Masonry arch bridge defect classes, with the underlying problems and potential solutions



algorithm. The image textures have also been resized to ensure a constant resolution in all the images. This resolution has been determined by standardising the average number of pixels for a brick course in each image. The standardised resolution used is 155 pixels per brick course, determined by the lowest resolution image in the dataset.

The image textures have been annotated with the different defect classes. This has been done by manually annotating the pixels in the images where a defect is present. An example of this for one of the annotated images is shown in Figure 2. The pixels containing mortar joints have also been annotated in order to use this dataset to determine the effect of mortar joints on defect detection performance.

Figure 2 Image annotations for recording defect locations showing; a) original image, b) mortar joints, c) crack locations, d) spalling locations, e) mortar loss locations, f) vegetation locations



### 2.3 Image Window Classification

The generated images have been segmented into smaller image patches, each 100 pixels in size, using a sliding window technique. The 100-pixel image size ensures that some image windows contain purely brick regions, while others contain a mixture of brick and mortar regions, as the height of each brick course in the images is 155 pixels. Each image window is assigned a class based on the annotations of the pixels it contains. Example image window patches for the different defect classes are shown in Figure 3.

The generated image windows are used to train a classifier to learn the different defect classes. The classifier used is a CNN which has been shown to perform well for classification tasks. The structure of CNNs have been inspired by the visual cortex, with the convolution layers of the model acting as feature extractors, simplifying the pixels of the input image into features which are then used to classify the image (Wang and Raj, 2017). The GoogleNet Inception v3 architecture (Szegedy et al., 2016) is used as it is one of the best performing models against the ImageNet classification benchmark. Only much more computationally expensive models have achieved slightly better performance (Canziani et al., 2017). This publicly available model has been pre-trained using the 1000 classes and 1.4 million images of the ImageNet dataset. Transfer learning is used on this pre-trained dataset as it means that a much smaller dataset can be used for training than would be necessary for training from scratch. Transfer learning fine tunes the pre-trained parameters based on the new classes and dataset. In this way much of the learning from the pre-training of the model can be applied to the new task of identifying

defects in masonry images. For training, 7000 image window patches for each defect class have been used.

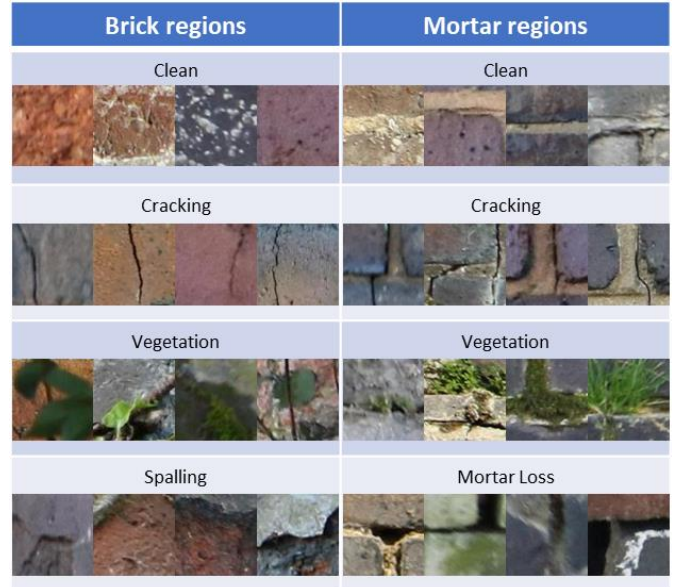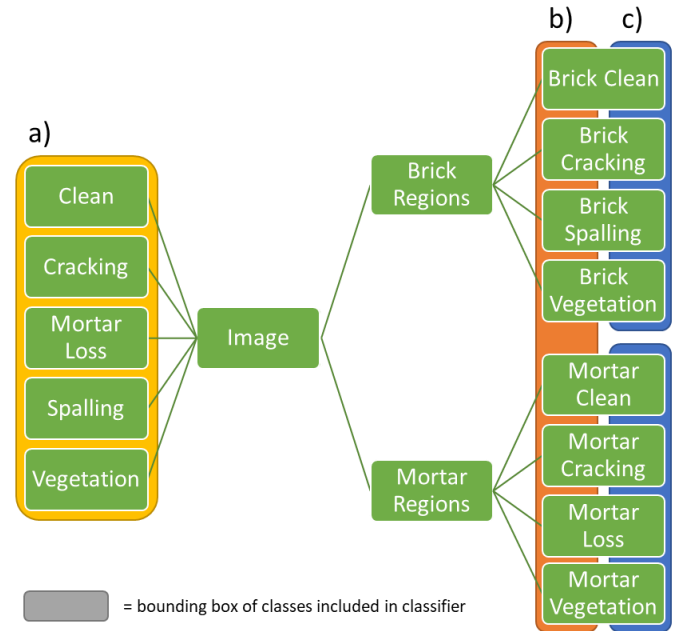Figure 3 Example image window patches for different defect classes



Figure 4 Different defect detection strategies used: a) no mortar/brick separation, b) mortar and brick defects labelled as separate categories, c) mortar regions and brick regions processed separately and merged after classification



### 2.4 Defect detection strategies

In order to determine the effect of mortar joints on defect detection accuracy, three different classification methodologies have been tested. These methodologies are summarised in Figure 4. The first strategy, shown in Figure 4(a), doesn't use any mortar joint information at all. Here only

the five different classes of defect are trained, with image window patches showing both mortar regions and brick regions being trained as the same relevant defect class. The second strategy is shown in Figure 4(b). Here separate defect classes are defined for mortar and brick regions, so that there are two defect classes for each defect type, one for the defect occurring in mortar and one for the defect occurring in brick. The final strategy, shown in Figure 4(c), completely separates the mortar and brick regions and uses a separate classifier for each. The two sets of classified images are then merged so that the image windows being classified are the same for all three detection methodologies.

For all three defect classification methodologies, only the image window patches that show fully brick regions or are centred on mortar regions are examined. Therefore, those image window patches that partly contain both brick and mortar regions are removed. As there is an overlap between image window patches, the whole of the masonry surface is still included 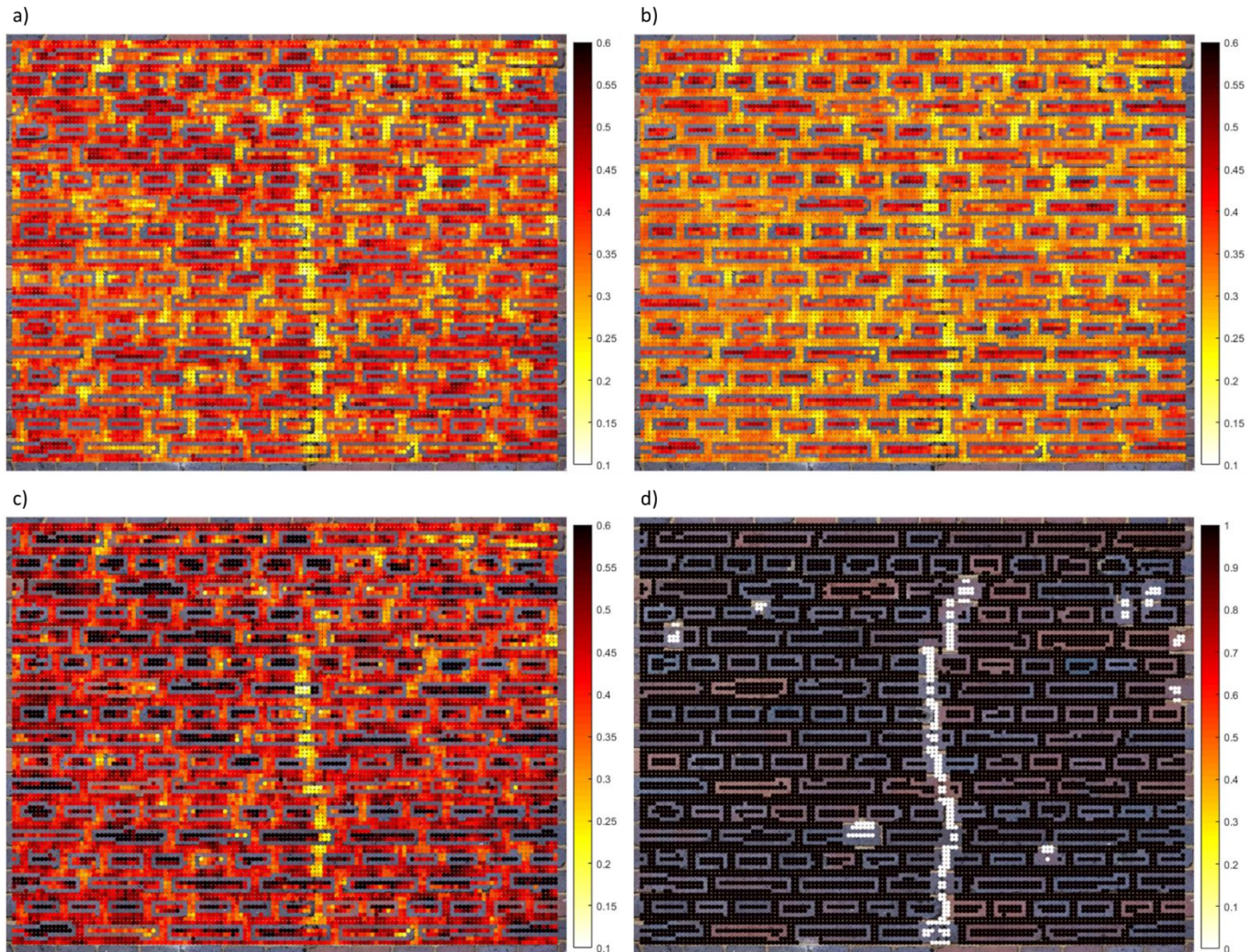in at least one examined patch. This step has been taken in order to ensure consistency of the data being examined by the three classification methodologies.

## 3. Results

The 94 annotated images have all been processed by the described methodology and their image window patches have been classified using the three classification strategies. For every image patch. the classifier assigns the probability that it belongs to each class of defect. This is then compared against ground truth data in order to determine the accuracy of classification for the three different classification strategies.

Figure 5 shows an example of the predicted output for the three different classification strategies, taken from one of the test images. To produce these plots, each image window has been assigned a shade based on the confidence that there are no defects, i.e. that it is a classified as a clean image window. The shaded image windows are plotted onto the test image at the centre point of the image window patch. It is therefore possible

Figure 5 Example output showing confidence of clean classification for different strategies: a) no mortar/brick separation, b) mortar and brick defects labelled as separate categories, c) mortar and brick regions processed separately and merged after classification, d) ground truth
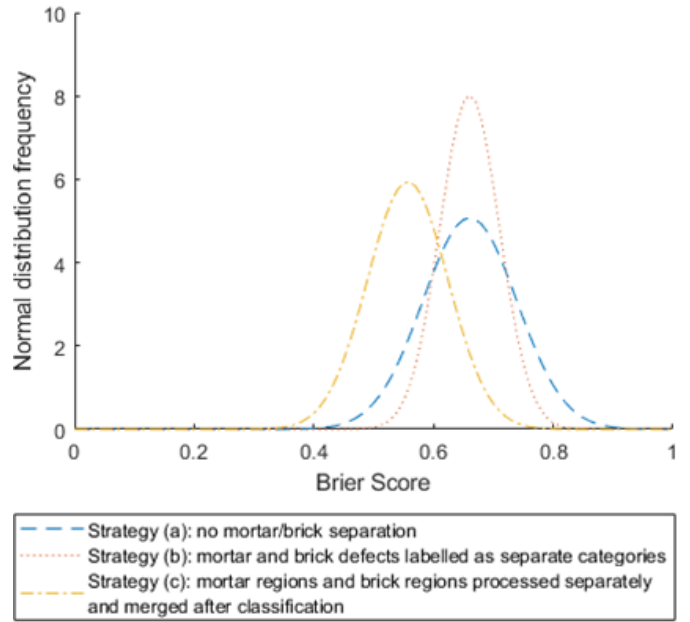


5

to visualise the regions of the test image that the image window patches refer to. The performance of the different classification strategies can be visualised by comparing Figure 5 (a - c) to Figure 5(d), the ground truth data. The ground truth data shows that there is a defect that runs down the length of the masonry image about in the centre. This is shown in all three of the outputs from the classifiers as a lighter area, meaning they have predicted a lower probability that the image windows are clean in this area. Contrasting Figure 5 (a - c), it is apparent that there is a larger contrast between the clean areas and the defect areas for the classification strategy where the mortar and brick regions are processed separately and then merged, than for the other two strategies. This is caused by a greater degree of confidence in the clean image windows being clean for this detection strategy. Additionally, for all three defect classification outputs, the brick areas are generally shaded darker than the mortar areas, meaning that they are predicted as more likely to be clean. This suggests that all three classification strategies are confusing the mortar areas with a class of defect.

This degree of confidence in predicting the correct category is measured by the Brier score. The Brier score measures the mean squared error between the predicted probability and the ground truth, for each defect class assigned in each image window. Its formulation for multi-category scoring is shown in Equation (1), where $p_{ic}$ is the predicted probability and $o_{ic}$ is the ground truth probability, for image window $i$ and class $c$. Here the total number of image windows and total number of classes are $N$ and $R$ respectively. For each image window, the ground truth probability for a defect class is assigned as one where the image window contains that class and zero where it does not.

$$Brier\ score = \frac{1}{N}\sum_{i=1}^{N}\sum_{c=1}^{R}(p_{ic} - o_{ic})^2 \qquad (1)$$

The Brier score, in this formulation takes values of between two as the worst score achievable and zero as the best score achievable. Figure 6 shows the Brier score that has been calculated for the three different defect classification strategies. Here, a slightly better Brier score is achieved by strategy c, (where the mortar regions and brick regions have been classified separately and merged after classification), then has been achieved by the other two strategies. Additionally, the distribution peak is higher, particularly for strategy b (where defects in mortar and brick regions have been labelled as separate classes), but also for strategy c, when compared to strategy a (where there is no mortar/brick separation). The higher peak is caused by a lower variance in the Brier score between the different images of bridges examined. This suggests that those classification strategies that incorporate mortar joint information (strategies b and c) are more consistent in performance, suggesting they cope better where the masonry images are noisier.

Figure 6 Brier score for different classification strategies



- - - Strategy (a): no mortar/brick separation
········· Strategy (b): mortar and brick defects labelled as separate categories
-·-·- Strategy (c): mortar regions and brick regions processed separately and merged after classification

Additionally, the performance of the three different classification strategies for correctly classifying the clean image window patches has been measured. Here, the predicted class for each image window is set as the class for which the predicted probability is the highest. Precision and Recall are measures of the performance of a binary classification. Precision (Equation 2) measures the proportion of the predicted instances of a class that are correctly predicted and recall (Equation 3) measures the proportion of the instances of a class that have been predicted. The $F_1$ score combines precision and recall as a measure of a classification's accuracy. It is computed as the harmonic mean of precision and recall (Equation 4). It takes values between zero at its worst and one at its best.

$$Precision = \frac{tp}{tp + fp} \qquad (2)$$
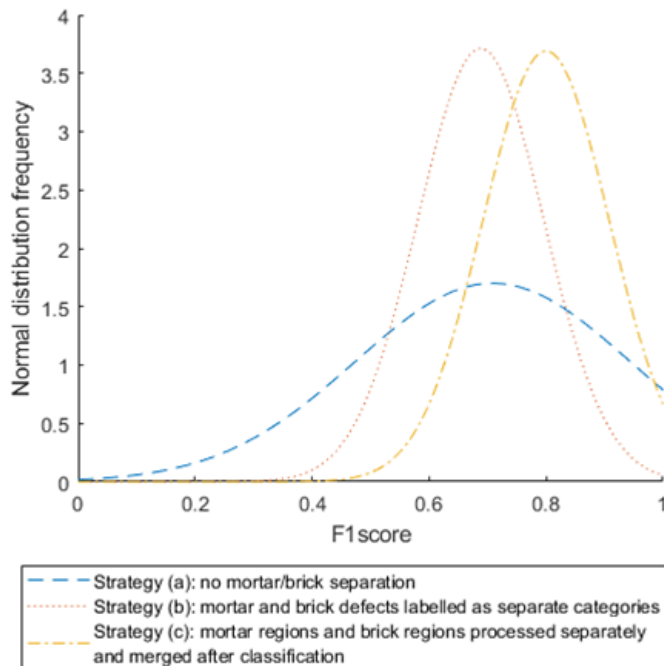
$$Recall = \frac{tp}{tp + fn} \qquad (3)$$

$$tp = true\ positive\ predictions$$
$$fp = false\ positive\ predictions$$
$$fn = false\ negative\ predictions$$

$$F_1\ score = 2 \times \frac{precision \times recall}{precision + recall} \qquad (4)$$

Figure 7 shows the $F_1$ score that has been computed for each of the three classification strategies for classifying clean image regions across all 94 test images. This suggests that strategy c (where mortar and brick regions have been classified separately) has a better performance than the other two defect classification strategies. As was the case for the Brier score, the results here also suggest that strategy a (where there is no

mortar/brick separation) is less consistent in its performance, due to the shorter and wider normal distribution profile.

Figure 7 $F_1$ score for different classification strategies for classifying clean image windows



## 4. Conclusions

This paper has initially reviewed the effect of the different types of observable defects on the structural condition of masonry arch bridges. This has been used to propose a workflow that can be used by a future automated bridge monitoring system to determine faults in a bridge and suggest appropriate remedial action based on a set of detectable symptoms. By using the proposed workflow, the main classes of defects in masonry that an automated visual detection system for masonry should be capable of detecting have been identified and have been used for the training of a CNN.

Three different defect detection strategies for separating the mortar and brick regions of masonry during classification have been used to determine the effect of the mortar joints on the performance of defect classification in masonry. Results suggest that separating the mortar and brick regions prior to classification causes an improvement in the confidence with which a classifier predicts masonry areas are clean. This leads to improvements in the Brier score and $F_1$ score for the classification. Additionally, less variation in the classification performance between different masonry images is found where the mortar and brick regions have been separated prior to classification, suggesting that this prior segmentation leads to the classifier performing better with noisier masonry images.

## 5. References

Canziani A *et al.* (2017) An Analysis of Deep Neural Network Models for Practical Applications. *2017 IEEE International Symposium on Circuits & Systems*.

Cha Y-J *et al.* (2017) Deep Learning-Based Crack Damage Detection Using Convolutional Neural Networks. *Computer-Aided Civil and Infrastructure Engineering* 32 (5): 361–78.

Chaiyasarn K *et al.* (2018) Crack Detection in Masonry Structures Using Convolutional Neural Networks and Support Vector Machines. *35th International Symposium on Automation and Robotics in Construction*.

Helmerich R *et al.* (2007) Guideliene for Inspection and Comndition Assessment of Existing European Railway Bridges. *Sustainable Bridges*, 1–259.

Koch C *et al.* (2015). A Review on Computer Vision Based Defect Detection and Condition Assessment of Concrete and Asphalt Civil Infrastructure. *Advanced Engineering Informatics* 29: 196–210.

Krizhevsky A *et al.* (2012). ImageNet Classification with Deep Convolutional Neural Networks. *25th International Conference on Neural Information Processing Systems*, 1097–1105.

Laefer DF *et al.* (2010). Reliability of Crack Detection Methods for Baseline Condition Assessments. *Journal of Infrastructure Systems* 16 (2): 129–37.

McCormick NJ *et al.* (2014). Assessing the Condition of Railway Assets Using DIFCAM: Results from Tunnel Examinations. *6th IET Conference on Railway Condition Monitoring (RCM 2014)*.

McKibbins L *et al.* (2006). *Masonry Arch Bridges: Condition Appraisal and Remedial Treatment (C656)*. CIRIA, London, UK.

McRobbie S (2008). Automated Inspection of Highway Structures. *Transport Research Laboratory*.

Network Rail (2017). *NR/L3/CIV/006. Handbook for the Examination of Structures*. UK.

Phares BM *et al.* (2004). Routine Highway Bridge Inspection Condition Documentation Accuracy and Reliability. *Journal of Bridge Engineering* 9 (4): 403–13.

Samy MP *et al.* (2016). Automatic Optical & Laser-Based Defect Detection and Classification in Brick Masonry Walls. *2016 IEEE Region 10 Conference (TENCON)*, 3521–24.

Szegedy C *et al.* (2016). Rethinking the Inception Architecture for Computer Vision. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2818–26. IEEE.

Wang H and Raj B (2017). On the Origin of Deep Learning. *ArXiv Preprint*.

Zhang L *et al.* (2016). Road Crack Detection Using Deep Convolutional Neural Network. *2016 IEEE International Conference on Image Processing (ICIP)*, 3708–12. IEEE.