

Folded Corner Detection Component

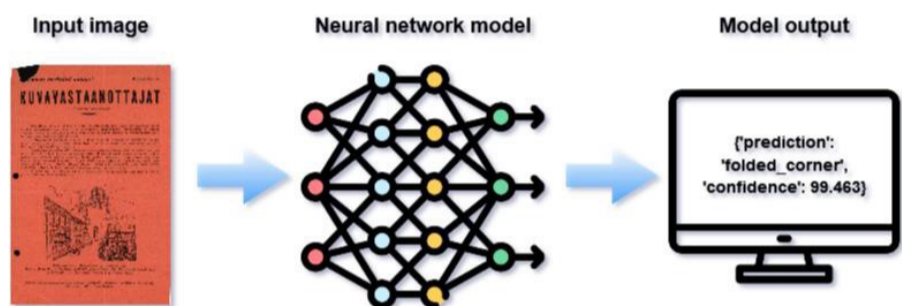
Background, benefits and purpose of use of the component

Documents in which the corner or edge of the paper is folded or torn may be scanned during the document digitisation process. If the fold or tear affects the readability of the content of the document, it is useful if such cases can be identified and, if necessary, scanned again. Machine identification of folded corners thus helps improve the quality of digital material and reduces the need for manual quality control.

What does the component do?

The component processes the image files received as input one at a time and forms a prediction of the content of the image based on the parameters it has learned during the training. If the component identifies a fold or tear in the image with a probability of more than 50%, the image is classified as incorrect. The component returns the predicted class to the user.

How is the component used?



Icons designed by Freepik and xnimrod from Flaticon

The component is available both as part of the Arkkiivi user interface (<http://www.arkkiivi.fi/>) and as a standalone application offering an application programming interface (API) on the DALAI project's GitHub page (<https://github.com/DALAI-project/CornerAPI>). More detailed instructions on how to use the component in different environments can be found on the attached website.

Component training

Image files in .jpg format aligned to 224 x 224 pixels have been used for component training. A total of about 35,000 documents have been used in the training, of which approximately 5,000 are documents containing folds or tears. The training material has been selected so that the component would identify folded and torn corners of the documents, regardless of whether the errors directly affect the readability of the document's content (e.g. a folded corner need not cover the text content).

Although efforts have been made to compile a variety of example cases for the training material, the amount of the material is limited and errors also occur in component classification. Among others, a document corner resembling a fold due to its colour or shape may lead to an incorrect classification.