

# Document Localization in Images Taken by Smartphones Using a Fully Convolutional Neural Network

Shima Baniadamdizaj (✉ [baniadam.shima@gmail.com](mailto:baniadam.shima@gmail.com))

Mohammadreza Soheili

Azadeh Mansouri

---

## Research Article

**Keywords:** Document corner localization, Smartphone Image Capturing, Deep Learning, Image Processing

**Posted Date:** October 4th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-952656/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# Document Localization in Images Taken by Smartphones Using a Fully Convolutional Neural Network

Shima Baniadam Dizaj

Department of electrical and computer engineering

Kharazmi University

Tehran, Iran

baniadam.shima@gmail.com

Mohammadreza Soheili

Department of electrical and computer engineering

Kharazmi University

Tehran, Iran

soheili@khu.ac.ir

Azadeh Mansouri

Department of electrical and computer engineering

Kharazmi University

Tehran, Iran

a\_mansouri@khu.ac.ir

**Abstract**— Today integration of facts from virtual and paper files may be very vital for the expertise control of efficient. This calls for the record to be localized at the photograph. Several strategies had been proposed to resolve this trouble; however, they may be primarily based totally on conventional photograph processing strategies that aren't sturdy to intense viewpoints and backgrounds. Deep Convolutional Neural Networks (CNNs), on the opposite hand, have demonstrated to be extraordinarily sturdy to versions in history and viewing attitude for item detection and classification responsibilities. We endorse new utilization of Neural Networks (NNs) for the localization trouble as a localization trouble. The proposed technique ought to even localize photos that don't have a very square shape. Also, we used a newly accrued dataset that has extra tough responsibilities internal and is in the direction of a slipshod user. The end result knowledgeable in 3 exclusive classes of photos and our proposed technique has 83% on average. The end result is as compared with the maximum famous record localization strategies and cell applications.

**Keywords**— Document corner localization, Smartphone Image Capturing, Deep Learning, Image Processing

## I. INTRODUCTION

Based on the benefit of use and portability of smartphones, enhancing the processing power, and enhancing the first-rate of pix taken the usage of smartphones, those telephones were capable of partly update the paintings of scanners in file imaging. In the meantime, because of the extraordinary talents of smartphones and scanners, there are troubles and demanding situations alongside the manner of turning telephones into scanners. Also, scanners are slow, expensive, and now no longer portable. Smart-telephones, on the opposite hand, have emerge as very clean to get to, use, or understand.

There are a few demanding situations to digitalize a file the usage of a phone, a number of which might be talked about. Lack of uniform light, shadows at the file, which may be the shadow of the hand, phone, or different objects, type of materials, colors and functions of the file, difference/extraordinary model withinside the heritage of the file and its contents, having three-D distortion, blurring, heritage complicated trouble of files along with coated pages, chessboard, etc., low file contrast, or bad phone digital digicam

first-rate, undetectable life of the file from its heritage because of being the identical color, light, etc., the complicated trouble of the file, for example, having folds, taking pics of multi-web page files along with books and identification cards, being a part of the file out of the picture, covered. Being a part of the file with the aid of using different objects, etc. A best approach ought to be sturdy and dependable for those demanding situations. It also can carry out on a phone in one of these affordable periods.

In standard, withinside the discipline of placing right into a laptop file, a few researchers were withinside the discipline of resolving or supporting to enhance the picture first-rate and lowering the troubles talked about/stated withinside the preceding paragraph, and others have given algorithms that during case of troubles withinside the picture taken with the aid of using the careless consumer, the file can nevertheless be located withinside the picture. There is a 3rd class of studies that, whilst enhancing picture first-rate and guiding the consumer to report with the aid of using a digital digicam or laptop the very best first-rate picture of the file, gives the set of laptop commands had to discover the file withinside the picture, that's a mixture of the preceding techniques.

There are extraordinary sorts of files. Digitalized paper files are simpler to carry, read, and proportion, search, index, and keep. One of the benefits of file imaging is the capacity to transform a paper file right into a virtual file, which with the aid of using persevering with the method of placing right into a laptop and the usage of virtual letter reputation techniques, the textual content and contents of the file may be effortlessly edited and searched. It is likewise viable to keep the virtual file withinside the area of outside reminiscence and effortlessly proportion or flow from one area to some other it among humans with a smaller overall area occupied with the aid of using something than the distance that a paper file occupies withinside the fitness of the Earth/the encircling conditions.

We make a contribution with the aid of using explaining algorithms that use system gaining knowledge of (ML) to gather information from scanned files on this paper. We additionally undergo a few realistic techniques that may be carried out the usage of software program assets like Python, PyTorch, TensorFlow, and OpenCV. This now no longer

simply highlights the strengths of system gaining knowledge of and deep gaining knowledge of on this situation, however additionally indicates that it's far a precious task. We endorse a way that makes use of deep convolutional neural networks to localize file corners in pix taken with the aid of using smartphones. Our approach outperforms the first-rate consequences amongst existed algorithms. It is greater standard in comparison to preceding techniques withinside the feel that it could be modified to match be sturdy and dependable to greater demanding situations with the aid of using schooling on higher consultant information.

## II. LITERATURE REVIEW

### A. Existed Datasets

To localize files in photos taken with the aid of using smartphones we want a real-international dataset this is accrued from a normal user. There are 4 exclusive datasets in record photos taken with the aid of using smartphones task. Three of those datasets include pics which have the identical or very near photos together. The fourth dataset accrued extra photos than others and additionally toward the real-international taken photos with exclusive demanding situations.

The to be had records set become used for the having to do with identifying the excellent of factors without measuring them with numbers method of identifying the worth, amount, or excellent of something of photos of files fascinated by smartphones [1]. The records set of Kumar et al. carries 29 exclusive files below exclusive angles and with blurring, and finally, 375 photos have been gotten. The dataset provided withinside the paper [2] makes use of 3 not unusual place sorts of paper in casting off exclusive sorts of distortion or harm which includes blurring, shaking, exclusive lights conditions, combining sorts of distortion in a picture, and taking photos which have one or extra distortions on the identical time, and the usage of extra than two, however now no longer a whole lot of sorts of smartphones, which makes these records set extra reliable.

Paper [3] gives a Mobile Identification Document Video records set (MIDV-500) made out of 500 movies for fifty wonderful identification record classes with floor truth, permitting examine on a huge type of record evaluation issues. The paper provides capabilities of the records set and method of identifying the worth, amount, or excellent of something outcomes for current techniques of face recognition, textual content line recognition, and records extraction from record fields. Because the sensitivity of identity papers is a vital aspect, as they encompass non-public records, all pix of supply files utilized in MIDV-500 are both now no longer copyrighted; unfastened to apply with the aid of using anyone, or disbursed below public copyright licenses.

The records set is provided withinside the article [4], which covers a few elements of the scene, which includes the lights conditions. A easy history become used. A robot arm become used to taking photos to get rid of the digital digicam shake. In the identical idea [5] provided a video dataset that There are five classes from easy to complicated, all with the identical content material and history, it carries films with 20 frames. And photos are pulled out or taken from something else from

those frames. Different smartphones have been used for the harm due to the device, and additionally with the aid of using the usage of exclusive files. A general of 4,260 exclusive photos of 30 files have been taken.

In the paper [6] a brand-new record dataset is provided this is toward the real-international photos taken with the aid of using users. The records separated and classified into easy, middle, and complicated obligations for detection. It carries nearly all demanding situations and carries exclusive record sizes and brands and backgrounds. It compares the end result of the record localizing techniques with famous techniques and cell laptop programs.

A newly posted dataset [7] this is the first Brazilian identification files public. All records withinside the BID Dataset is from faux records to conform with the non-public statistics privateness legislation. This attempt intends to boost up studies development in identity record picture processing with the aid of using permitting teachers to freely make use of the BID Dataset of their examine.

In the paper [6] a new document dataset is presented that is closer to the real-world images taken by users. The data separated and labeled into simple, middle, and complex tasks for detection. It contains almost all challenges and contains different document sizes and types and backgrounds. It compares the result of the document localizing methods with well-known methods and mobile computer programs.

A newly published dataset [7] that is the first Brazilian identification documents public. All data in the BID Dataset is from fake data to comply with the personal information privacy legislation. This effort intends to accelerate research progress in identification document image processing by allowing academics to freely utilize the BID Dataset in their study.

### B. Algorithms and methods

Due to the demanding situations, it isn't feasible to digitize files the usage of smartphones without preprocessing or post-processing, and count on desirable consequences in all situations. That is why algorithms were proposed to enhance the consequences. The impact of picture graph venture algorithms at the end result may be divided into 3 categories: 1. Reduce demanding situations earlier than taking pictures 2. Fixed problems at the same time as taking pics 3. Solve demanding situations after taking pictures.

One of the earliest strategies of report localization changed into primarily based totally on a version of the heritage for segmentation. The heritage changed into modeled via way of means of taking an photo of the heritage without the report. The assessment among each pix changed into used to decide wherein the paper changed into found. This approach had the plain negative aspects that the digital digicam needed to be saved desk bound and pix needed to be taken. [8]

In general, the algorithms used to locate the report withinside the photo may be divided into 3 categories: 1. use of extra hardware 2. depend on photo processing strategies 3. take benefit of gadget mastering strategies. This trouble has arisen

with the unfold of smartphones from 2002 to 2020 and may be improved..

### 1) *Additional Hardware*

In the article [9] they gift courses for the person in taking with fewer demanding situations primarily based totally on specific functions. As a result, the photo calls for plenty much less pre-processing to localize the document. This approach became now no longer very person-pleasant for the customers because of the constraints and slowdown of digitization. Article [10] used this approach for localizing. After pre-processing, extra algorithms are required to finish the localization task. These algorithms are categorized as follows: 1. Use of more hardware 2. Make use of device imaginative and prescient techniques 3. The utility of deep mastering algorithms.

A scanning utility is presented [11] that consists of actual-time web page recognition, fine assessment, and automated detection of a web page cover [12] whilst scanning books. Additionally, a transportable tool for putting the cellular smartphone all through scanning is presented. Another paper that used extra hardware introduces a scale-invariant function remodel into the paper detection gadget. [13] The hardware of the paper detection gadget includes a virtual sign processor and a complicated programmable common-sense tool. The tool is able to obtaining and processing images. The software program of this gadget makes use of the SIFT approach to discover the papers. Compared to the conventional approach, this set of rules offers higher with the detection process. In the paper [14] paper detection desires a sheet of paper with a few styles revealed on it. It makes laptop imaginative and prescient one step in the direction of getting used withinside the actual world.

### 2) *Machine vision techniques*

The set of rules [15] works with the aid of using finding capability line segments from horizontal experiment strains. Recognized line segments are multiplied or blended with textual content line segments from adjoining experiment strains to shape large textual content blocks, which might be then subjected to filtering and refinement. The paper [16] introduces a popularity gadget in complicated historical past video pics.

The proposed morphological approach [17] is insensitive to noise, skew, and textual content orientation. It is likewise freed from artifacts which are typically delivered with the aid of using each a constant/fine international thresholding manner and a block-primarily based totally constant-length nearby thresholding manner. A morphology-primarily based totally approach is proposed [18] to extract essential assessment capabilities as clues for locating the preferred license plates. The assessment function is touchy to modifications in illumination and invariant to diverse variations consisting of scaling, translation, and skew. The paper [19] applies part detection and makes use of a low threshold cost to clear out non-textual content edges. Then a nearby threshold is selected to each preserve low-assessment textual content and to simplify the complicated historical past of the excessive-assessment textual content. Next, textual content vicinity enhancement operators are proposed to focus on the one's regions with both excessive part thickness or excessive part density.

A step-with the aid of using-step approach is defined [20] wherein candidate regions are searched from the enter picture the usage of gradient records after which the plate vicinity most of the applicants is decided and the boundary of the vicinity is adjusted with the aid of using placing a plate template. In the textual content excerpts from the video picture paper [21] the vertices of the chosen video pics are recognized. After deleting a few remoted corners, the last corners are joined collectively to shape textual content candidate regions. Target pics [22] are decided on at constant periods from pics detected with the aid of using a scene extrude detection approach. For every decided-on body, a satiation histogram is used to phase with the aid of using satiation grouping across the satiation histogram.

The approach [23] locates candidate areas immediately withinside the DCT-compressed area the usage of the depth version records encoded withinside the DCT area. The work [24] makes use of a clean picture historical past to discover areas of interest (RoI). [25] proposes a linear length phase detector that achieves dependable consequences with a low quantity of fake detections and calls for no parameter tuning. This set of rules may be examined and in comparison, with modern-day algorithms on a huge set of herbal pics. An method to figuring out file applicants [26], defined in a given picture, the usage of Geodesic Object Proposals [27]. The enter pics were down-sampled to extract exciting structures/capabilities, lessen noise, and growth runtime pace and accuracy. The consequences confirmed that it's miles promising to apply Geodesic Object Proposals withinside the assignment of file item popularity. Furthermore, [28] operators are associated with the max-tree and min-tree representations of files in pics. An easy approach for computing the tree of paperwork is defined in article [29]; it really works for nD pics and has a quasi-linear complexity whilst statistics quantization is low.

The methodology [30] is primarily based totally on projection profiles in mixture with a linked thing marking process. Signal cross-correlation is likewise used to confirm the detected noisy textual content regions. Several extraordinary steps are used for assignment [31]: a preprocessing manner the usage of a low-byskip Wiener filter, a difficult estimation of the foreground regions, a calculation of the historical past vicinity with the aid of using interpolation of adjoining historical past intensities, a threshold calculation with the aid of using combining the calculated historical past floor with the authentic picture together with an picture up-sampling and subsequently a post-processing step to enhance the pleasant of the textual content areas and to preserve the road connectivity. It has been proposed that every horizontal textual content line crosses a specific set of vertical strains at non-horizontal locations to keep away from the skew difficulty in virtual files [32]. They create a correlation matrix and calculate the file's skew attitude with excessive precision the usage of handiest the pixels on vertical strains. For the whisker files [33] the authors advanced a sturdy function-primarily based totally approach to mechanically merge a couple of overlapping pics. The intended set of rules [34] makes use of a mixture of combinatorial meeting of capability quadrangle applicants from a group of line segments and projective file

reconstruction with a regarded body rate. Fast Hough Transform [35] is used for line popularity. A 1D amendment of the threshold detector is proposed for the set of rules. Three localization algorithms are defined within the paper [36].

All algorithms employ key points, and of them moreover look at nearly horizontal and nearly vertical strains within the picture. The counseled approach [37] is a way for detecting the 4 corners of a file in sensible scenario the usage of a coarse-to-high-quality localization strategy. In the primary stage, the 4 corners are kind of anticipated with the aid of using a deep neural community primarily based totally on Joint Corner Detector (JCD) with an interest mechanism that kind of localizes the file place through an interest map. In [38], the authors advocate an set of rules that approaches an picture in conditions wherein modern-day quadrilateral file border detection algorithms aren't optimized for instances wherein one of the file borders is absolutely out of the body, blurred or low assessment. Also, display that this method decreases the danger of wrong detections with the aid of using 34% at the MIDV-50 dataset.

CREASE, a content-conscious file rectification method that optimizes a per-pixel attitude regression loss, a curvature estimation loss, and a 3-d coordinate estimation loss for offering picture rectification maps, became delivered with the aid of using Markovitz and et.al [39]. Their counseled approach complements OCR accuracy, geometric, and visible similarity-primarily based totally metrics with the aid of using repairing folded and creased files with clues from the file's nearby and international scale characteristics. The first step of the proposed version is used to expect 3-d structure, angles, and curvature, at the same time as the second one section is used to expect the backward map. They use a pixel-degree attitude regression loss, that is beneficial in each 3-d estimation and cease-to-cess training. The 3-d estimation approach additionally learns the attitude side-assignment at the phrases within the textual content, optimizing for clarity within the rectified picture, whilst the curvature evaluation side-assignment suits the attitude regression with the aid of using mapping its divergences.

### 3) Machine Learning method

The have a look at [40] gives a CNN-primarily based totally method for correct real-time report positioning and treats the trouble as a key factor identity task. The 4 corners of the files are anticipated collectively through a Deep Convolutional Neural Network. The paper [41] first identifies the sort of report and classifies the photos, and then, through understanding the sort of report, a coordinated localization approach for the report is performed, which allows the extraction of information.

Among the processes for automated report photo processing, algorithms primarily based totally on Convolutional Neural Networks (CNN) were extensively correctly applied. The first CNN method to report reputation become proposed through LeCun et al. [42]. Later, Kang et al. [43] supplied a CNN for report classification. Tensmeyer [44] supplied a CNN method to binarizing report photos. In every other have a look at this is beneath Neath publication, we used distinct convolutional neural networks for the identical task. In that work, we used DeepLabv3 [45] for segmentation and

localization files in photos taken through smartphones. This CNN includes 3 distinct parts. For the semantic segmentation part, it's far feasible to apply distinct CNNs and we determined to use MobileNetv2 because of its real-time results.

## III. METHODOLOGY

We model the problem of document localization as corner detection. The method needs a new ground truth model as a mask of the document part and the non-document part. We demonstrate the document with white (255) and non-document parts with black (0).

### A. Dataset Preprocessing

We use [1-5] datasets as train and validation datasets and make the image size similar using the max image height and width among images using zero paddings. And use [6] dataset as the test dataset for the evaluation and comparing the proposed method with the previous methods and mobile applications. During train, we use data augmentation by random scaling of input pictures (from 0.5 to 2.0) and randomized left-right scaling flipping. And also, by randomly rotating the images.

### B. Proposed method

For the undertaking of file localization in snap shots taken with the aid of using smartphones, we used U-Net [46] and the fine-tuning approach to retrain a few ultimate layers within the U-Net deep neural community. This community has benefited from deconvolution. In this dissertation, we've got taken into consideration locating the placement of the file within the snap shots as a semantic segmentation undertaking. Convolutional neural networks tested affordable accuracy within the segmentation of semantic photograph data.

Figure 1 indicates the community structure. It is created of a contracting route (at the left) and an expansive route (at the right). The convolutional community's forming path follows the usual structure. It is split into  $3 \times 3$  convolutions (unpadded) which can be implemented repeatedly, every accompanied with the aid of using a rectified linear unit (ReLU) and a  $2 \times 2$  max-pooling step with stride 2 for down-sampling. We growth the wide variety of characteristic channels with every down-sampling segment. So, every step within the increasing course starts off evolved with an up-sampling of the characteristic map, following with the aid of using a  $2 \times 2$  convolution ("up-convolution") that cuts the wide variety of characteristic channels in 1/2 of, and a concatenation with  $3 \times 3$  convolutions. every observed with the aid of using a ReLU, and the for this reason cropped characteristic map from the contracting route Due to the lack of nook pixels in any convolution, cropping is needed. A eleven convolution is hired within the ultimate layer to map every 64-thing characteristic map to the specified wide variety of classes. The community consists of a complete of 23 convolutional layers.

Because of structure and pixel-primarily based totally photograph segmentation constructed from convolutional neural community layers, U-Net outperforms different traditional models. It works additionally with snap shots from

a small dataset. The layout of this structure started out with the exam of scientific photographs. The approach of lowering the peak and width dimensions is properly-known. The dimensionality discount segment withinside the top and width that we follow withinside the convolutional neural community - that is, the pooling layer - is applied withinside the 2d 1/2 of of the version withinside the shape of a size rise, as is properly known. By retaining the wide variety of channels withinside the enter matrix constant, the pooling layer decreases top and width data. The calculation is a level withinside the manner of lowering complexity. Each member of the photo matrix is called a pixel. In brief, a pixel that displays businesses of pixels is called the pooling layer.

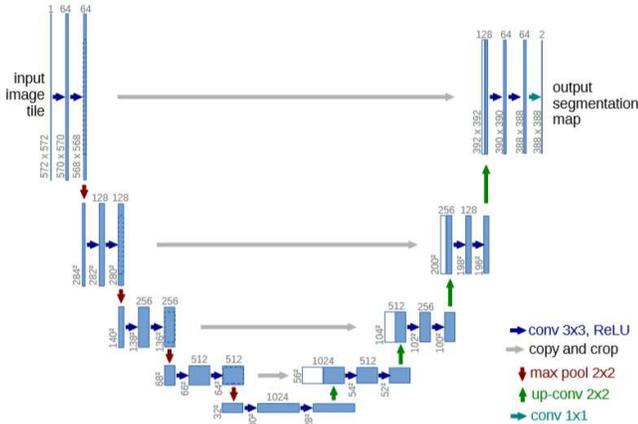


Figure 1 shows a U-net design (example for 32x32 pixels at the lowest resolution). A multi-channel feature map is represented by each blue box. The number of channels is indicated on the box's top. The x-y size is supplied at the box's lower-left border. The white box indicates feature maps that have been replicated. The arrows denote the different operations. [44]

These layers are designed to grow the output decision. For localization, the sampled output is blended with high-decision functions withinside the model. A sequential convolution layer attempts to generate a extra specific output primarily based totally in this knowledge. U-Net receives its call from the structure, which resembles the letter U while visualized, as visible withinside the diagram above. A segmented output map is made from the enter photos. In the second one half, the maximum particular characteristic of the structure. There is not any absolutely connected layer withinside the network. The convolution layers are the handiest ones which are included. A ReLU activation characteristic turns on every traditional convolution process.

The length of the enter statistics is reduced first in a conventional autoencoder design, observed with the aid of using the subsequent layers. The structure's encoder segment is now complete, and the decoder component begins. This element teaches linear characteristic representation, with the scale growing progressively. Enation, with the scale step by step increasing. The output length equals the enter length on the quit of the structure. This structure is best for preserving output length, however it has one flaw: it linearly compresses the statistics, ensuing in a hassle in which all functions aren't transmitted-Net plays deconvolution at the decoder aspect (i.e., withinside the 2nd half) and, besides, can conquer this bottleneck problem, As a end result, functions are misplaced

through connections from the encoder aspect of the design. It is likewise hard to outline limitations as components of the identical elegance rub up towards every other. It is suggested to apply values with a huge weight withinside the loss characteristic for that reason while first isolating the statistics to be segmented from the context. As a end result, U-Net can localize photos which have 3-D distortion. So, the very last end result isn't a square form region.

### 1) Contracting/Down-Sampling layer

The convolutional network's contracting route follows the conventional architecture. It includes 3x3 convolutions (unpadded convolutions) which might be implemented repeatedly, every operation is observed with the aid of using a rectified linear unit (ReLU) and a 2x2 max pooling process with stride of two for down sampling. We double the quantity of characteristic channels with every down-sampling phase. This contracting route pursuits to seize the enter image's heritage in order that segmentation may be performed.

By doubling the scale of filters in every step, this route is used to extract extra functions of the image. While getting into the subsequent degree we are able to do 2x2 max-pooling to get the most pixel value, for this reason dropping a few functions, however maintaining the most pixel value. So, on the remaining layer of Down-sampling, we have become the lower-degree functions of an image.

### 2) Expansive/Up-Sampling layer

It, too, is made of a couple of enlargement blocks, so just like the contraction layer. After up sampling the characteristic map, a 2x2 convolution (upward convolution) is carried out, which halves the variety of characteristic channels, observed via way of means of a concatenation with the enter photo, that's flippantly cropped from the course of contraction. and 3x3 convolutions, every observed via way of means of a ReLU withinside the expansive course. The intention of this increasing course is to permit particular localization while evaluating facts from the contracting course.

We might be doing up-sampling withinside the equal level to hold the equal length because the photo. We additionally have the pixel feature values for all the agencies after down-sampling. There's no want to be worried due to the fact we have got ignored any capabilities in down-sampling via way of means of the use of max-pooling. Up-sampling restores the whole photograph via way of means of means of copying the feature map of a stage with the equal down-sampling filters to the extent with the equal up-sampling filters, retaining the capabilities. As a result, we get the whole photograph again and may pinpoint wherein the disorder is withinside the photo for every class. Transpose convolution is the time period for this. Convolution, on the alternative hand, is used to examine the sizable photo. Therefore, while up sampling, every characteristic layer on a discounted sampling facet is introduced to the corresponding characteristic layer at the top sampling facet so one can achieve the entire decision photo and therefore localize the class. Implementing place capabilities with the photo's international information.

### C. Training Protocol

The education stage wishes a one-of-a-kind floor fact from the paper [6] so we offer a masked floor fact (Fig.2) that the report element and the non-report element Are differentiated with black and white colors. the report with white (255) and the non-report element with black (0). After freezing the supposed layers, the very last community has been up to date and applied in Ubuntu Linux model 16.04 LTS implementation and programming. The STRIX-GTX1080-O8G pix card and the Core i7 6900k processor with 32 GB of RAM also are used for education, power, and testing.

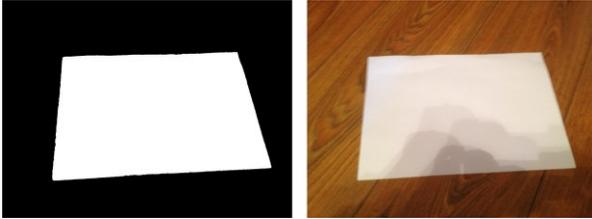


Fig 1. Image sample with masked ground truth

## IV. DISCUSSION AND RESULTS

### A. Evaluation Protocol

We implement a test to evaluate the system's quality on a test set, unseen contexts, and unseen documents. The IoU method described in [47], has been used to evaluate. First, the perspective effect is deleted from the ground truth (G) and predicted (S) with the help of image size. We call new situations (G') and (S') respectively so that the IoU or Jaccard index is equal to:

$$\text{IoU} = (\text{area}(G' \cap S')) / (\text{area}(G' \cup S')) \quad (1)$$

The final result is the average of the IoU value for each image.

### B. Results

We run the version on a take a look at dataset and examine our effects to the formerly launched effects at the equal dataset [6]. While this test is enough to assess the device's accuracy at the opposition dataset, it does have drawbacks. One, it can't display us how powerful our device deflates to unknown contents due to the fact numerous samples of contents had been used for schooling. Second, it's far not able to offer facts approximately how nicely our framework generalizes to unseen files with comparable content material. We cross-validate our approach via way of means of deleting every content material from the schooling set after which checking at the deleted content material to restoration the weaknesses and degree generalization on unseen content material.

The end result is as compared with all famous methods, algorithms, and cellular programs that could resolve the report localization assignment in snap shots. The different methods' effects are as compared primarily based totally on [6] (Fig.3). Our approach effectively generalizes to unseen files, as effects are proven in Fig.4. This isn't always unreasonable for the reason that the low decision of the enter photo prevents our version from counting on functions of the report's format or content material. We additionally finish from the effects that

our approach generalizes nicely to unseen easy contents. It is vital to say that the approach is designed to be powerful on a low wide variety of sources like a midrange phone with out the use of cloud or server-primarily based totally recourses. The frames from the 4 corners are processed withinside the order withinside the implementation. It can enforce snap shots suddenly via way of means of strolling a batch of 4 snap shots thru the version. This ought to bring about a enormous growth in inefficiency.

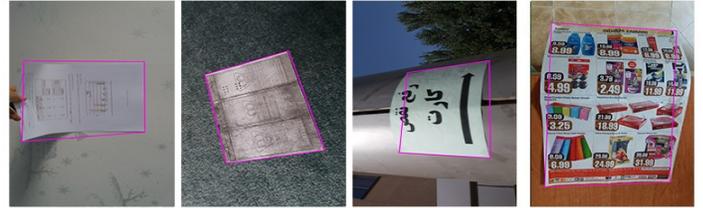


Fig 2. The result of the proposed method on dataset [6]

In Table 1, the final results in different categories are presented and Fig.4 shows the result in comparison with the previous methods. And also, there are different samples from different categories of results using the proposed method in Fig.3.

Table 1. The result of the proposed method on dataset [6]

Simple	Moderate	Difficult	Average
100%	84%	65%	83%

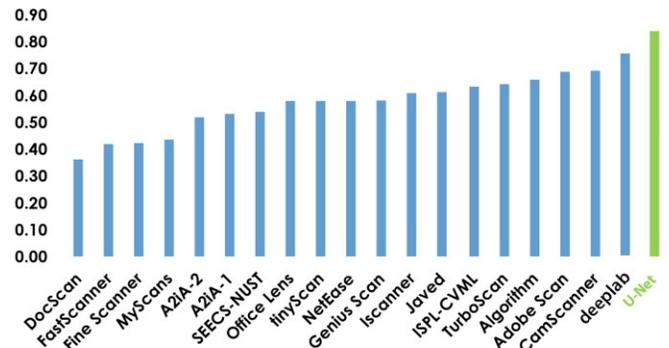


Fig 4. The result of the proposed method compared with previous methods

## V. CONCLUSION AND FUTURE WORKS

In this paper, we offered a brand-new software of U-Net for file localization in pictures taken through smartphones. The very last end result is the great end result of localization strategies. We make contributions through explaining algorithms that use gadget learning (ML) to find statistics from scanned files. We additionally undergo a few realistic strategies that may be implemented the use of software program assets like Python, PyTorch, TensorFlow, and OpenCV. We used all dependable datasets on this task. And, the primary dataset to examine is the newly accrued dataset with a diverse range of file localization challenges. Also, we gift a unique technique for finding files in herbal pictures. The hassle of localization is modeled as a hassle of key factor detecting. We display that this technique could make assumptions properly on new and unseen files through the use of a deep convolutional network.

Moreover, we illustrated that generalization on formerly unseen complicated structures is possible. Besides this, we expect that enhancing generalization on unseen complicated files may be performed through both the use of extra complicated pictures withinside the education or artificially destroying the education pictures through the use of patches of diverse colors. The very last end result of the evaluation takes a look at set suggests that the proposed technique has 100% curacy on easy pictures and the accuracy on common is 89% this is the great end result amongst stated outcomes till now. To boom the overall performance of our set of rules withinside the future, we can recall the geometrical functions of the textual content and content.

#### REFERENCES

- [1] Sheikh, Hamid R., Muhammad F. Sabir, and Alan C. Bovik. "A statistical evaluation of recent full reference image quality assessment algorithms." *IEEE Transactions on image processing* 15.11 (2006): 3440-3451.
- [2] Ye, Peng, and David Doermann. "Document image quality assessment: A brief survey." 2013 12th International Conference on Document Analysis and Recognition. IEEE, 2013.
- [3] Arlazarov, Vladimir Viktorovich, et al. "MIDV-500: a dataset for identity document analysis and recognition on mobile devices in video stream." *Компьютерная оптика* 43.5 (2019).
- [4] Nayef, Nibal, et al. "SmartDoc-QA: A dataset for quality assessment of smartphone captured document images-single and multiple distortions." 2015 13th International Conference on Document Analysis and Recognition (ICDAR). IEEE, 2015.
- [5] JC. Burie, J. Chazalon, et al. "ICDAR2015 competition on smartphone document capture and OCR (SmartDoc)." 13th International Conference on Document Analysis and Recognition, IEEE, 2015.
- [6] Dizaj, Shima Baniadam, Mohammadreza Soheili, and Azadeh Mansouri. "A New Image Dataset for Document Corner Localization." 2020 International Conference on Machine Vision and Image Processing (MVIP). IEEE, 2020.
- [7] de Sá Soares, Álysson, Ricardo Batista das Neves Junior, and Byron Leite Dantas Bezerra. "BID Dataset: a challenge dataset for document processing tasks." *Anais Estendidos do XXXIII Conference on Graphics, Patterns and Images*. SBC, 2020.
- [8] H. Lampert H, T. Braun, et al. "Oblivious document capture and real-time retrieval." *Proceedings. Camera-Based Document Analysis and Recognition*, 79-86, 2005.
- [9] Chen, Francine, et al. "SmartDCap: semi-automatic capture of higher quality document images from a smartphone." *Proceedings of the 2013 international conference on Intelligent user interfaces*. ACM, 2013.
- [10] Jayaraman, Dinesh, et al. "Objective quality assessment of multiply distorted images." 2012 Conference record of the forty sixth asilomar conference on signals, systems and computers (ASILOMAR). IEEE, 2012.
- [11] Kleber, Florian, et al. "Mass Digitization of Archival Documents using Mobile Phones." *Proceedings of the 4th International Workshop on Historical Document Imaging and Processing*. ACM, 2017.
- [12] Fototechnischer Ausschuss der KLA. 2016. *Wirtschaftliche Digitalisierung in Archiven*. (2016).
- [13] Zhu, Juan, Shuai Wang, and Fanying Meng. "SIFT method for paper detection system." *Multimedia Technology (ICMT)*, 2011 International Conference on. IEEE, 2011.
- [14] Quan, Neng, Xiaolong Zhou, and Shengyong Chen. "Scan paperback books by a camera." *Information and Automation (ICIA)*, 2016 IEEE International Conference on. IEEE, 2016.
- [15] Zunino, Rodolfo, and Stefano Rovetta. "Vector quantization for license-plate location and image coding." *IEEE Transactions on Industrial Electronics* 47.1 (2000): 159-167.
- [16] Kuwano, Hidetaka, et al. "Telop-on-demand: Video structuring and retrieval based on text recognition." 2000 IEEE International Conference on Multimedia and Expo. ICME2000. *Proceedings. Latest Advances in the Fast-Changing World of Multimedia (Cat. No. 00TH8532)*. Vol. 2. IEEE, 2000.
- [17] Hasan, Yassin MY, and Lina J. Karam. "Morphological text extraction from images." *IEEE Transactions on Image Processing* 9.11 (2000): 1978-1983.
- [18] Hsieh, Jun-Wei, Shih-Hao Yu, and Yung-Sheng Chen. "Morphology-based license plate detection from complex scenes." *Object recognition supported by user interaction for service robots*. Vol. 3. IEEE, 2002.
- [19] Cai, Min, Jiqiang Song, and Michael R. Lyu. "A new approach for video text detection." *Proceedings. International Conference on Image Processing*. Vol. 1. IEEE, 2002.
- [20] Kim, Sunghoon, et al. "A robust license-plate extraction method under complex image conditions." *Object recognition supported by user interaction for service robots*. Vol. 3. IEEE, 2002.
- [21] Hua, Xian-Sheng, et al. "Automatic location of text in video frames." *Proceedings of the 2001 ACM workshops on Multimedia: multimedia information retrieval*. ACM, 2001.
- [22] Kim, Hae-Kwang. "Efficient automatic text location method and content-based indexing and structuring of video database." *Journal of Visual Communication and Image Representation* 7.4 (1996): 336-344.
- [23] Zhong, Yu, Hongjiang Zhang, and Anil K. Jain. "Automatic caption localization in compressed video." *IEEE transactions on pattern analysis and machine intelligence* 22.4 (2000): 385-392.
- [24] Wu, Victor, R. Manmatha, and Edward M. Riseman. "Finding text in images." *ACM DL*. 1997.
- [25] Von Gioi, Rafael Grompone, et al. "LSD: A fast line segment detector with a false detection control." *IEEE transactions on pattern analysis and machine intelligence* 32.4 (2010): 722-732.
- [26] Leal, Luciano RS, and Byron LD Bezerra. "Smartphone camera document detection via Geodesic Object Proposals." *Computational Intelligence (LA-CCI)*, 2016 IEEE Latin American Conference on. IEEE, 2016.
- [27] Krähenbühl, Philipp, and Vladlen Koltun. "Geodesic object proposals." *European conference on computer vision*. Springer, Cham, 2014.
- [28] Carlinet, Edwin, and Thierry Géraud. "A comparative review of component tree computation algorithms." *IEEE Transactions on Image Processing* 23.9 (2014):3885-3895.
- [29] Géraud, Thierry, et al. "A quasi-linear algorithm to compute the tree of shapes of nD images." *International symposium on mathematical morphology and its applications to signal and image processing*. Springer, Berlin, Heidelberg, 2013.
- [30] Stamatopoulos, N., B. Gatos, and A. Kesidis. "Automatic borders detection of camera document images." 2nd International Workshop on Camera-Based Document Analysis and Recognition, Curitiba, Brazil. 2007.
- [31] Gatos, Basilios, Ioannis Pratikakis, and Stavros J. Perantonis. "Adaptive degraded document image binarization." *Pattern recognition* 39.3 (2006): 317-327.
- [32] Chang, Fu, Chun-Jen Chen, and Chi-Jen Lu. "A linear-time component-labeling algorithm using contour tracing technique." *computer vision and image understanding* 93.2 (2004): 206-220.
- [33] Zhang, Zhengyou, and Li-Wei He. "Whiteboard scanning and image enhancement." (2016).
- [34] Skoryukina, Natalya, et al. "Real time rectangular document detection on mobile devices." *Seventh International Conference on Machine Vision (ICMV 2014)*. Vol.9445. International Society for Optics and Photonics, 2015.
- [35] Duda, Richard O., and Peter E. Hart. "Use of the Hough transformation to detect lines and curves in pictures." *Communications of the ACM* 15.1 (1972): 11-15.
- [36] Skoryukina, Natalya, et al. "Document localization algorithms based on feature points and straight lines." *Tenth International Conference on Machine Vision (ICMV2017)*. Vol. 10696. International Society for Optics and Photonics, 2018.

- [37] Zhu, Anna, et al. "Coarse-to-fine document localization in natural scene image with regional attention and recursive corner refinement." *International Journal on Document Analysis and Recognition (IJ DAR)* 22.3 (2019): 351-360.
- [38] Tropin, D. V., Konovalenko, I. A., Skoryukina, N. S., Nikolaev, D. P., & Arlazarov, V. V. (2021, January). Improved algorithm of ID card detection by a priori knowledge of the document aspect ratio. In *Thirteenth International Conference on Machine Vision (Vol. 11605, p. 116051F)*. International Society for Optics and Photonics.
- [39] Markovitz, A., Lavi, I., Perel, O., Mazor, S., & Litman, R. (2020, August). Can You Read Me Now? Content Aware Rectification Using Angle Supervision. In *European Conference on Computer Vision* (pp. 208-223). Springer, Cham.
- [40] Javed, Khurram, and Faisal Shafait. "Real-Time Document Localization in Natural Images by Recursive Application of a CNN." *Document Analysis and Recognition (ICDAR), 2017 14th IAPR International Conference on*. Vol. 1. IEEE, 2017.
- [41] Awal, Ahmad Montaser, et al. "Complex document classification and localization application on identity document images." *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. Vol. 1. IEEE, 2017.
- [42] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [43] L. Kang, J. Kumar, P. Ye, Y. Li, and D. Doermann, "Convolutional neural networks for document image classification," in *2014 22nd International Conference on Pattern Recognition*. IEEE, 2014, pp. 3168–3172.
- [44] C. Tensmeyer and T. Martinez, "Document image binarization with fully convolutional neural networks," in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 1.
- [45] Chen, Liang-Chieh, et al. "Rethinking atrous convolution for semantic image segmentation." *arXiv preprint arXiv:1706.05587* (2017).
- [46] Ronneberger, Olaf, Philipp Fischer, and Thomas Brox. "U-net: Convolutional networks for biomedical image segmentation." *International Conference on Medical image computing and computer-assisted intervention*. Springer, Cham, 2015.
- [47] Rezatofghi, Hamid, et al. "Generalized intersection over union: A metric and a loss for bounding box regression." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019.