

Text Recognition - Real World Data and Where to Find Them

Klara Janouskova, Jiri Matas

Centre for Machine Perception, Department of Cybernetics
Czech Technical University, Prague, Czech Republic

Lluís Gomez, Dimosthenis Karatzas

Computer Vision Center
Universitat Autònoma de Barcelona, Spain

Abstract—We present a method for exploiting weakly annotated images to improve text extraction pipelines. The approach uses an arbitrary end-to-end text recognition system to obtain text region proposals and their, possibly erroneous, transcriptions. The proposed method includes matching of imprecise transcription to weak annotations and edit distance guided neighbourhood search. It produces nearly error-free, localised instances of scene text, which we treat as “pseudo ground truth” (PGT).

We apply the method to two weakly-annotated datasets. Training with the extracted PGT consistently improves the accuracy of a state of the art recognition model, by 3.7 % on average, across different benchmark datasets (image domains) and 24.5 % on one of the weakly annotated datasets.

I. INTRODUCTION

Written text is an important source of information for humans and plays an important role in everyday life, being frequently present in scenes with man-made structures. It is one of the most common classes in general object detection datasets like CoCo [1]. Research in end-to-end reading systems is a very active research field in academia and industry alike. It is an essential part of many applications ranging from translation systems and autonomous driving to image retrieval or visual question answering.

In recent years, with the introduction of deep neural network models, the field has advanced substantially. At the same time, state of the art performance comes at the cost of the requirement of large-scale annotated data for training. Such data need to be rich in geometry, style and content. Typical ground truth data is defined at the granularity of words as polygonal regions in the image along with the corresponding text transcriptions. The acquisition of such data requires substantial human effort and is very costly.

The lack of data is usually approached in two different ways, either by generating synthetic data as in [2], [3], [4], [5], [6], [7] or with different forms of weakly, semi or unsupervised learning on real data as in [8], [9], [10].

While fully annotated real world data are expensive and sparse, weakly annotated data in the form of images along with a set of words likely to appear in them are common. An example of a source for such weakly annotated data are product databases where we can readily obtain the name of the product and other meta-data. A weakly annotated dataset, based on a product database, used in our experiments is shown in Figure 1, other sources in Figure 2. Images from mapping services like Google Maps are another potential source where

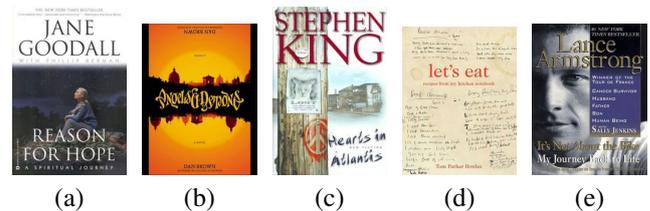
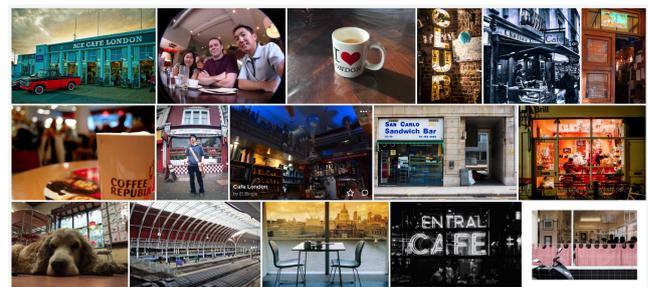


Fig. 1: The ABC Dataset images are diverse, some have (a) both a simple layout and font, (b) a very artistic, almost illegible font and (c) a font that resembles handwriting. Other have (d) both dense background and hand-written text or (e) the background varies significantly even at the word level. The dataset includes weak annotations - the name of the author and the title. Not all text is annotated, location of the annotated text is not provided.



(a) Searching for ‘Cafe London’ photos in the Flickr application.



(b) Artist and album name in Amazon Music product database.

Fig. 2: Potential sources of weakly annotated data; the availability of such data for research purposes changes over time.

street names and numbers or business names are words very likely to appear in an image and easily obtainable through location-based search. For example, given images from the location where a restaurant is supposed to be, it is expected that the name of the restaurant will be visible in some of them

In this paper, we present a new method for automatically generating pseudo ground truth (PGT) in the form of text regions and their transcription from weakly annotated data consisting of text transcriptions only with no information about the text location. The weak annotations may be noisy

- incomplete (not all the text in the image is there) and with distractor words (text that is not in the image). The method requires an end-to-end text recognition system (E2E) pre-trained with annotated data. In contrary to our method which uses automatically obtained weak annotations, previous methods used weak labels which alleviate but do not eliminate human participation, such as annotating the areas of interest.

The core idea of our method is that the OCR output of the recognition model can be used to identify the most probable text match from the weak annotations by finding the one with the lowest edit distance. The detections that produce an exact match with the weak label are assumed to be correct. Furthermore, the recognition output can be used to find the modifications of the detected regions that minimize the edit distance to the matched text. For example, if we have predicted the word ‘car’ and the best match was ‘cartoon’, running the recognition again on the detected region extended to the right may decrease the distance between the matched and predicted text, possibly leading to the prediction of the matched word ‘cartoon’. We have experimentally verified that the probability of the recognition model giving the same output as the weak label for a wrong text region is very low, thus it is safe to use such regions as PGT for further training.

In summary, given an image and a dictionary of words as an input, the method outputs a subset of the dictionary words with their corresponding text regions in the image. The method is independent of the underlying implementation of the models.

Possible applications of our method are improving the performance of an existing general E2E system or domain adaptation, where the source domain has full ground truth data whereas only weak labels are available for the target domain.

We apply our method to data from two different sources - a database of images of book-covers downloaded from Amazon books, using the title and the author of the book as weak annotations, and the Uber-Text dataset [11], which is very similar to the kind of data that could be obtained using a mapping system. We train a recognition model with the PGT generated from both sources on various benchmarks, showing that it consistently improves the recognition accuracy across a wide range of datasets. We chose Uber-Text, a dataset with region-level annotations available, to evaluate how our approach compares to fully supervised training in ablation studies. Otherwise, the full annotations are not used.

The contributions of the paper are:

- We present a novel method for automatic generation of pseudo ground truth data from images with weak annotations.
- The accuracy of the best-performing open-source model is improved significantly (3.7 % on average) on a wide range of benchmark datasets by training with the PGT generated from two datasets.
- The PGT data significantly improve recognition performance in weakly-supervised domain adaptation (24.5 %).
- The PGT annotations in the Amazon Book Covers dataset will be published.

II. RELATED WORK

We first give a short introduction of scene text detection and recognition methods, and an overview of methods for generating synthetic data. Then we focus on weakly-supervised learning for text detection and recognition.

A. Text Detection and Recognition

Before the deep learning era, methods based on SWT or MSERs were used for text detection, for example [12], [13]. Subsequent models were mostly based on region proposal approaches like [14]. Recently, methods have rather turned to segmentation-based approaches like [15] and focused on representing arbitrarily shaped text, for example [8], [16].

Recent approaches for text recognition rely on deep learning. Most methods can be described by 4 stages. Transformation - a Spatial Transformer Network [17] normalizing the input image. Feature extraction - a CNN such as VGG [18] or ResNet [19] maps the input image to feature maps. Sequence modelling - BiLSTMs are used to provide contextual information to the feature maps. Prediction - either CTC [20] or attention-based prediction [21] is used to convert the encoded features into a character sequence. Some methods treat the two tasks jointly, sharing features for both detection and recognition, for example [22], [23], [24].

B. Synthetic data for text detection and recognition

The work of [5] and [14] had a great influence on the performance of text detection and recognition systems. Synthetic data have proven to be very effective for training generic text localisation systems. Still, the lack of realism (both in terms of positioning, and blending with the scene), diversity (in terms of text styles and scene backgrounds) and contextualisation of the text in the scene, have been limiting factors. More recent works aim to improve some of these aspects [2], [4], [6], [7], or exploit instead real scene text data to do augmentation [3], still, cannot replace the quality of real-world data.

C. Weakly supervised learning

The proposed method builds on top of pseudo-labelling techniques [25], a simple strategy for semi-supervised learning where part of the data is fully labelled and Pseudo-Labels are created for unlabelled data as the class with the maximum predicted probability and further treated as true labels.

Focusing on Chinese street view images, [26] have rough location and transcription of some text instances annotated. An online proposal matching module is incorporated in the whole model. The main difference from our method is that they do not do any modification of the proposed regions.

In [24], an existing OCR engine different from the one being trained is used to provide partial labels for one million unlabelled images. The partially labelled data is then used to train the recognition part of an E2E model, improving the results significantly. The method relies on a confidence threshold to filter out noisy labels while our method relies on weak annotations to minimize the risk of incorrect labelling.

Focusing on text detection, [10] propose multiple approaches for unsupervised and weakly supervised learning. Their unsupervised approach simply relies on filtering out predictions with low confidence score. An improved approach requires weak annotations, where regions containing text are annotated and it is known that regions distant from the annotated ones do not contain any text, allowing for more accurate false positives filtering. Their last approach relies on rectangular bounding boxes as weak labels.

III. DATASETS

We introduce the datasets used in our experiments.

Amazon Book Covers (ABC) is a dataset created from more than 200,000 images downloaded from Amazon Books. The author and the title of each book serve as weak annotations. The same data were already used for genre prediction in [27]. Some illustration images are shown in Figure 1.

Uber-Text dataset (UT) is one of the biggest datasets for text detection and recognition. It contains 117,969 images with 571,534 labelled text instances split into training, validation and test sets. Each set is divided into two subsets according to the image resolution - either 1K or 4K. The images were obtained through the Bing Maps Streetside program and come from 6 different cities in the US. The annotations are line-level. Most of the text regions form semantic units such as business names, street signs or street numbers. The datasets contains a lot of not annotated text, some text regions are not annotated at all, while some readable text is labeled as unreadable [11].

MJSynth (MJ) contains almost 9M synthetically generated images of English words for text recognition. The text generation process performs the following steps: Font rendering, border/shadow rendering, coloring, projective distortion, natural data blending and noise introduction [14]. **SynthText (ST)** is a synthetic dataset designed for scene-text detection, widely used for recognition, too. It has over 7M text instances in 8,000 images [5]. **Synthetic Multi-Language in Natural Scene Dataset (MLT)** contains 245,000 images in total with text instances in multiple scripts: Arabic, Bangla, Chinese, Japanese, Korean and Latin. The dataset was published in [23] and the authors have adapted the framework of [5]. A non-latin dictionary was used and it contains special, non-alpha-numeric characters. We only use the Latin script subset of the dataset, which contains 288,917 text instances in total [23].

IIIT 5K-word (IIIT) is a collection of 5,000 cropped words from Google image search using queries which are likely to contain text [28]. The training set consists of 2,000 images, the remaining 3,000 form the test set. **Street View Text (SVT)** was collected from the Google Street View, providing annotators with a lexicon for each image, containing texts such as business names. Only the words from the lexicon were localised and provided with transcription, the rest of the text is ignored. There are 257 and 647 images of cropped words in the training and test sets [29]. **Street View Text - Perspective (SVT-P)** is a dataset of 645 images collected from Google Street View focused on perspective projections [30].

ICDAR2003 (IC03), collected for the ICDAR 2003 Robust Reading competitions [31], has 258 training and 251 testing images with 1,156 and 1,110 annotated words respectively. **ICDAR2013 (IC13)** is a dataset with 'focused text', the text being the main content of the image. It consists of a training set of 229 images with 848 words and a test set of 233 images with 1095 words. [32]. In contrast, **ICDAR2015 (IC15)**, focuses on incidental scene-text - the images were not taken with text in mind. The training set contains 1,000 images (4,468 words) and the test contains 500 images (2,077 words) [33].

Total-Text (TT) is a dataset of 1,555 scene images with 9,330 annotated words. The images were collected with curved text in mind and the images often contain texts of different orientations [34]. **CUTE80 (CT)** contains 80 images with 288 words, focusing on curved text [35].

IV. PGT GENERATION

This section describes the pseudo ground truth (PGT) generation algorithm (PGT-GEN). The algorithm uses weakly annotated images and an existing end-to-end reading system (E2E). All the steps are executed independently for each image. First, we define the E2E output and the structure of the weak annotations. Then we describe the algorithm and its components in detail.

Given an image I , the output $O = \{(bb_1, tt_1), \dots, (bb_t, tt_t)\}$ of the end-to-end reading system is a set of t text bounding box predictions and the corresponding text transcriptions. The transcriptions $T = (tt_1, \dots, tt_t)$ are strings (possibly containing spaces) and the bounding boxes are oriented rectangles. It is necessary that the recognition output from a bounding box bb can be obtained independently: $REC(I, bb) = tt$.

Each image is associated with a list of texts $A = (t_1, t_2 \dots t_n)$ where each text $t_i = (w_1, w_2, \dots, w_m)$ is a non-empty ordered sequence of words. The set of weak labels $G = \bigcup_{i=1}^n g_i$ is obtained as a union of sets of k -grams, $k \in \{1, \dots, 5\}$. Each set of k -grams g_i is formed by strings (consecutive words from t_i), sub-sequences of t_i of length k joined into a single string by the space character. In the simplest of cases, each text t_i only consists of a single word but because the texts are assumed to be extracted automatically as metadata accompanying the images, it may even be multiple words that form a semantic unit — a name of a product, its description, a business' name, contact information. These words are likely to appear in the image close to each other and get merged by the detector.

For example, the texts could be "Sherlock Holmes" and "221B Baker Street". A and G would then be

$$A = ((\text{"Sherlock"}, \text{"Holmes"}), (\text{"221B"}, \text{"Baker"}, \text{"Street"}))$$

$$G = \{\text{"Sherlock"}, \text{"Holmes"}, \text{"Sherlock Holmes"}, \text{"221B"}, \text{"Baker"}, \text{"Street"}, \text{"221B Baker"}, \text{"Baker Street"}, \text{"221B Baker Street"}\}.$$

A. PGT-GEN algorithm

The PGT-GEN algorithm takes the image I , E2E output O and the set of weak labels represented as k-grams G as an input and outputs the PGT - a localized subset of G .

Algorithm 1: PGT-GEN

Input: I, O, G
Output: PGT
 $P := \text{AssignWeak}(O, G);$
 $PGT = \{\};$
foreach $((bb, tt), g) \in P$ **do**
 $(bb_f, tt_f) = \text{FindOptimalBox}(I, bb, tt, g);$
 if $\text{IsPGT}(tt_f, g)$ **then**
 $PGT = PGT \cup \{(bb_f, g)\};$
 end
end
return PGT

AssignWeak - Weak annotation assignment. Each element from O is assigned at most one weak annotation from G . We construct a directed bipartite graph $B_G = (V, E)$ between O and G , thus $V = O \cup G$. For each output $o \in O$, $o = (bb, tt)$ and weak annotation $g \in G$ it holds that

$$(o, g) \in E \iff \text{dist}(tt, g) = \min_{i=1}^{|G|} \text{dist}(tt, g_i) \quad (1)$$

$$(g, o) \in E \iff \text{dist}(tt, g) = \min_{i=1}^{|T|} \text{dist}(tt_i, g) \quad (2)$$

where dist is the Levenshtein distance.

Then, a set of proposals P is created:

$$P = \bigcup_{i=1}^{|O|} \text{Assign}(o_i, E) \quad (3)$$

$$\text{Assign}(o, E) = \begin{cases} \emptyset & \text{for } W(o, E) = \emptyset \\ [W(o, E)]_R & \text{otherwise} \end{cases} \quad (4)$$

$$W(o, E) = \{(o, g) : (o, g) \in E \wedge (g, o) \in E \wedge \text{match}(o, g)\}. \quad (5)$$

We define $\text{match}((bb, tt), g) = \frac{\text{dist}(tt, g)}{\max(\text{len}(tt), \text{len}(g))} < 1$ to filter out completely irrelevant proposals and $[\cdot]_R$ selects an element from a set randomly. In most cases, $|W(o, E)| \in \{0, 1\}$.

At this point, we could apply some simple filtering to the set of proposals P instead of the edit distance guided neighbourhood search, for example, select

$$P' = \{p \in P, p = ((bb, tt), g) | \text{dist}(tt, g) = 0\} \quad (6)$$

and then output

$$PGT = \bigcup_{((bb, tt), g) \in P'} (bb, tt). \quad (7)$$

This would be equivalent to selecting proposals where the predicted transcription was equivalent to the weak label text



Fig. 3: Improved localization of weak labels by the neighbourhood search - the detected bounding boxes (blue) and the transformed ones (green).

for PGT - we implement this version and compare it to the proposed one, showing the superiority of the proposed method.

FindOptimalBox - Edit distance guided neighbourhood search. For each proposal $((bb, tt), g) \in P$, we search for an optimal bounding box bb_f which minimizes the Levenshtein distance between the recognized text tt_f and g .

If $\text{dist}(tt, g) = 0$, we assume that bb is already optimal and assign $bb_f = bb$. If not, we predefine a set of new boxes in the neighbourhood of bb and run the recognition on those in parallel, selecting one with minimal distance from g for bb_f . The generation of the set of predefined boxes is explained in detail in Appendix A.

We compute $tt_f = \text{REC}(I, bb_f)$ and the normalized edit distance between tt_f and g as $d = \frac{\text{dist}(tt_f, g)}{\max(\text{len}(tt_f), \text{len}(g))}$.

Finally, we find out whether (bb_f, tt_f) satisfies our requirements for being a PGT (**IsPGT**) as:

$$\text{IsPGT}(tt_f, g) = \begin{cases} \text{True} & \text{for } d = 0 \vee \text{isClose}(d, tt_f, g) \\ \text{False} & \text{otherwise} \end{cases} \quad (8)$$

where $\text{isClose}(d, tt_f, g) = (d < \theta \wedge |tt_f| > \lambda \wedge tt_f^0 = g^0 \wedge tt_f^{-1} = g^{-1})$, s^0, s^{-1} are the first and the last characters of a string s . The thresholds θ, λ are set to $\theta = 0.35, \lambda = 4$. The intuition behind the IsClose function is that even if the recognized text tt_f and the assigned text g are not identical, it is possible that there was simply an error in the recognition step. If the relative edit distance between two longer texts is low and the first and the last characters are the same, it is likely that tt_f should actually be g .

Examples of how the neighbourhood search aids the PGT generation process can be seen in Figure 3.

V. END-TO-END READING SYSTEM

In this section, the end-to-end reading system (E2E) used in our experiments is described. Separate models for detection and recognition are used.

A. Detection

For text detection, we adopt TextSnake [16]. It is based on a fully convolutional network - U-net with VGG-16 [18] backbone - which estimates the geometry attributes of text instances. A text instance is described as a sequence of ordered, overlapping disks centered at symmetric axes. Each disk is associated with a potentially variable radius and orientation.

The following values are predicted for each pixel: tcl, tr, r and α , corresponding to the text center line, text region, radius and angle. Thresholds varying on different datasets are applied to tcl and tr to obtain binary masks tcl_b and tr_b . Focusing

on straight text, we replace the proposed subsequent post-processing steps for text instance reconstruction with a method based on least squares fitting of the text center line points, which improves the orientation of the final bounding box.

First we obtain the connected components CC from $tr_b * tcl_b$. For each component $cc \in CC$, where cc are all the component points, we estimate the bounding box (c_x, c_y, w, h, α) directly.

The angle α is estimated by total least squares fitting of a line to the points of cc shifted by the mean value of the coordinates $m = (m_x, m_y)$ to the origin, using the slope of the best fitting line as the bounding box angle.

The height h is determined via the biggest radius predicted within the component: $h = \max(r[cc]) * 2$. To determine the width w , we project the shifted points onto the line with slope α passing through the origin and find the projected vectors with maximum norm in both directions, p_{pos} and p_{neg} . The width is calculated as $w = |p_{pos} - p_{neg}| + h$. The extra h is added to the width because during training, the tcl is shrank by $\frac{h}{2}$ (assuming the radius is constant, $\frac{h}{2}$, for straight text with rectangular ground truth). Afterwards, we shift p_{pos}, p_{neg} back by m and calculate the center of the bounding box as the middle point: $(c_x, c_y) = m + \frac{p_{pos} + p_{neg}}{2}$.

B. Recognition

We adopt the best performing architecture from [36], similar to STAR-net [37] but with a different prediction mechanism.

1) *Transformation*: An input image I is transformed into a rectified image \tilde{I} . It predicts the parameters of a thin-plate spline (TPS) transformation, a variant of spatial transformer network (STN) [17]. The whole module consists of a localization network, a grid generator and a grid sampler.

We use grayscale images as the input and both the input and output dimension are fixed to 32×150 pixels. For more details, we refer the reader to [38], [36], [37].

2) *Feature Extraction*: Given the rectified image \tilde{I} , the feature extractor outputs a feature map

$$V = \text{CNN}(\tilde{I}) = \{v_i\}, i = 1, \dots, K \quad (9)$$

where $K = 38$ is the number of columns in the output feature map (512×38).

3) *Sequence modeling*: a BiLSTM network [39] creates contextual features from the visual features v_i and outputs $H = \text{Seq}(V)$. We use a 2-layer BiLSTM. An i^{th} layer identifies two hidden states: forward $h_i^{(t),f}$ and backward $h_i^{(t),b}, \forall t$. A fully-connected layer between the two BiLSTM layers determines one hidden state, $\hat{h}_t^{(i)}$, from $h_i^{(t),f}$ and $h_i^{(t),b}$. The dimension of the hidden states and the FC layer is 256.

4) *Prediction*: Finally, a single layer LSTM [40] attention decoder produces the output sequence of characters $Y = y_1, y_2, \dots, y_n, y_t = \text{softmax}(W_o s_t + b_o)$, where W_o and b_o are trainable parameters, and s_t is the decoder LSTM hidden state at time t . The decoding stops when the (EOS) symbol is emitted. The model is trained with the cross entropy loss function. For more details on the attention mechanism, please see [8], [21], [24].

This recognition model is used in all of our experiments and we will refer to it simply as OCR.

VI. EXPERIMENTS

The pseudo ground truth (PGT) generation method was tested with two different sources of weakly annotated data, the Amazon book covers dataset (ABC) and the Uber-Text (UT) training set where we ignored localization information.

The detection part of E2E (TextSnake) was trained on a mix of SynthText, ICDAR2015 and Total-Text datasets. The post-processing thresholds of TextSnake were set to $tr = 0.4$ and $tcl = 0.7$, which lead to the best PGT generation performance on a small subset of UT and ABC images. We do not filter out the text marked as unreadable or don't care during training to maximize the use of available data. Detection recall is more important than precision for PGT generation – the more words detected, the more potential pseudo-labelled examples are available. On the other hand, false positives are very unlikely to be matched against weak annotations, thus they have minimal impact, besides slowing down the process.

The recognition part (OCR), which also serves as a base-line model (OCR_b) in our experiments, was trained on the ST (Synth-text), MJ (MjSynth) and MLT (Synthetic Multi-Language in Natural Scene) datasets. The OCR_b recognizes 70 characters – letters, not distinguishing lower and upper case, digits and frequent special characters like punctuation, brackets, the (EOS) symbol and the space.

Most recognition datasets provide word-level annotations, and thus space is never part of the transcription. We included space in the character set for three different reasons. First, if the model is capable of predicting spaces, it helps to guide the PGT generation process - a bounding box that is too wide leads to a space being predicted at the beginning or end of the transcription. Second, if the detector merges horizontally adjacent words, the recognizer often splits the text by recognizing a space between the merged words. Third, it allows exploiting annotations that contain multiple words.

Synthetic datasets used for training have word-level annotations and thus provide no training data for the space character. We therefore extended some of the bounding-boxes and included spaces at the beginning and end of the ground truth annotations. This produced a model with a limited ability to recognize the space. It was further improved during training on PGT, since it contains multi-word texts. During evaluation, we strip any leading/trailing spaces from the predictions.

The OCR processes images with a fixed resolution of 32×150 . Input images are first resized isotropically to height 32. If the width of the resized image is less than 150, the image is extended to the left and padded with zeros. If the width exceeds 150, it is horizontally shrunk to 150 - only in this case the aspect ratio of the input images is not preserved. The procedure of training with PGT is explained in Figure 4.

A. PGT from the Uber-Text dataset

The experiment evaluates the PGT method as an adaptation technique to the Uber-Text dataset domain. The performance

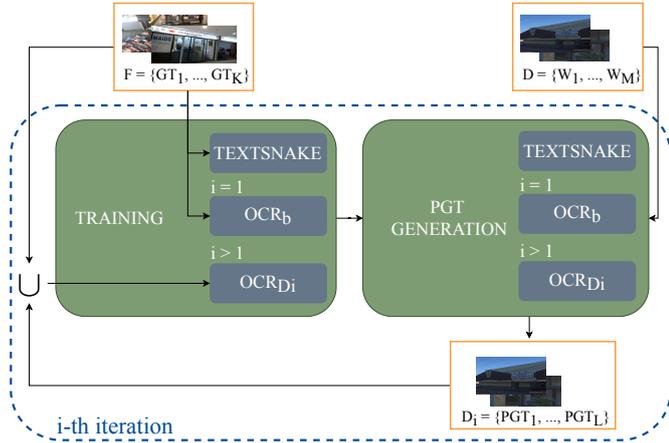


Fig. 4: Training procedure. In the i -th iteration, $i > 1$, the outputs of the $OCR_{D_{i-1}}$ and TextSnake models are used to generate a pseudo ground truth (PGT) dataset D_i from a weakly annotated dataset D . A new model OCR_{D_i} is trained on the union of D_i and fully annotated data F . In the first iteration, OCR_b and TextSnake are pretrained on F .

is also compared to fully-supervised and semi-supervised training in the UT domain.

A reference method, OCR_{UT_F} , is obtained by fully supervised training on UT_F , the set of 138,437 transcriptions and corresponding rectangular crops from the UT training set that contain no unreadable characters in transcription. The crops are the minimum area enclosing rectangles of the ground truth polygons. Appendix B contains details of a semi-supervised learning via pseudo-labelling experiment. The best performing model from this experiment is referred to as $OCR_{UT_{PL}}$. OCR_{UT_F} and $OCR_{UT_{PL}}$, as well as other OCRs described in this section, are validated on a set of 5,000 random transcriptions from the UT validation set.

For PGT generation, the whole UT training set is used. To facilitate GPU computations, we split large (about 4K) images into 16 blocks ensuring no text instance is split and discard those with no text. Such empty blocks are common since text instances are sparse in many images. Each of the original ground truth transcriptions is a weak label in our experiment, the ground truth polygons are discarded. The weak labels are transformed into a set of n-grams, as explained in the PGT generation section. N-grams containing the $*$ symbol (unreadable or unknown characters) are discarded.

PGT generation and OCR training. In the first iteration, 92,909 PGT text instances were obtained. The number of PGT text instances increased in all iterations, reaching 113,810 texts after six iterations when the OCR performance stopped improving – a summary is shown in Figure 5. The recognition rate, calculated on 20,000 randomly selected transcriptions from the UT test set, increased in each iteration from the baseline 41.6 % to 66.1 % in the sixth iteration. The accuracy of the fully supervised OCR_{UT_F} and semi-supervised $OCR_{UT_{PL}}$ is 78 % and 45.7 %, respectively. The PGT has reduced the gap between the baseline model and the fully-supervised one

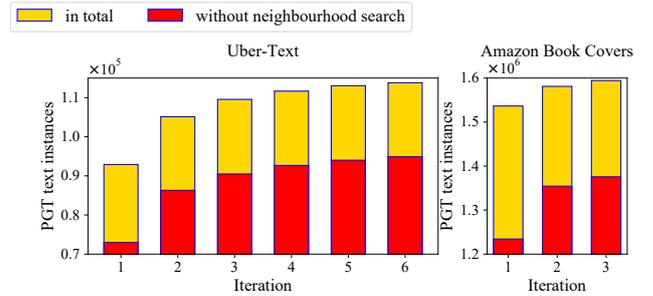


Fig. 5: Number of PGT text instances on the UT and the ABC datasets. On UT, the text location information is ignored.

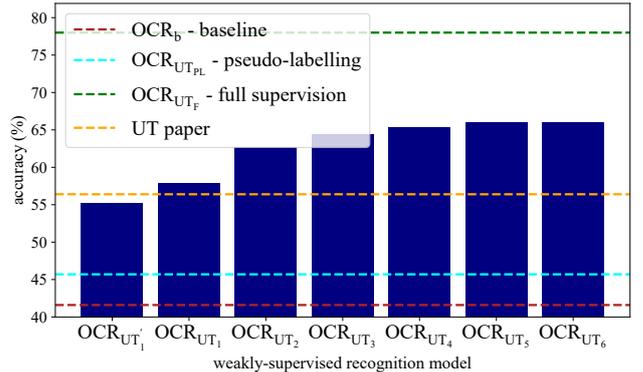


Fig. 6: Recognition rates on the Uber-Text test set. The model trained in iteration i is denoted as OCR_{UT_i} . UT'_1 is a subset of UT_1 obtained without the neighbourhood search. UT paper is the model of [11].

by 67 %. The performance of PGT training is limited by the detector which was not trained on the new domain. Improving the detector may help to reduce the gap further.

To test the contribution of the neighbourhood search and of allowing imperfect matches, a dataset, denoted UT'_1 , is created. It contains only the detections that matched with 0 edit distance with some of the weak labels. The accuracy of the model trained with UT'_1 is 2.7 % lower, showing the importance of the additional retrieved PGT text.

Figure 6 compares the accuracy of different models and Tables III and IV summarize the different experiments.

B. PGT from book covers

The PGT is generated from the whole ABC dataset - over 200,000 images, using the author and the title of each book as weak label. The title often includes a subtitle - while the author and title are almost always present in the image, the subtitle is less common, or there may only be a part of it.

We performed three iterations, again increasing the number of PGT generated in each of them. In the last one, ABC_3 , a dataset of 1,594,333 cropped images of text is generated. Figure 5 contains a summary of the PGT generation results and Table III contains a more detailed results.

C. PGT accuracy

We have analyzed the accuracy of the PGT method on UT_6 and ABC_3 , 500 text instance crops from each. We have found 21 and 10 images with a wrong PGT in the ABC and the UT datasets, respectively. Those are images that do not contain text or images where some characters/punctuation in the PGT are wrong. The majority of these errors are due to false positives/very blurred texts, leading to a prediction of a common word such as ‘the’, ‘on’ with low confidence and thus can be filtered out. For more details, see Appendix C.

D. Results on benchmark datasets

A recognition model trained with the PGT data from the previous experiments is evaluated on different domains using various commonly used recognition datasets - IIIT 5K-word (IIIT), Street View Text (SVT), ICDAR2003 (IC03), ICDAR2013 (IC13), ICDAR2015 (IC15), Street View Text - Perspective (SVT-P) and CUTE80 (CT). We evaluate on test set subsets commonly used by researchers as identified in [36]. All the models were validated on a union of the training sets of all the previously mentioned datasets. The reported metric is the percentage of correctly recognized words.

To evaluate models trained on word-level data only, a test set of 20,000 images where each image only contains a single word, referred to as UT_w , is also created.

We also evaluate on the UT_w subset of the UT dataset but it was not included in the validation set. We remove any spaces from all the predictions and when evaluating on datasets with images that contain punctuation but the ground truth does not, we filter any non-alphanumeric characters out.

Training with either UT_1 or ABC_1 generated in the first iterations of the PGT generation consistently improves the performance. On some datasets, UT boosts the accuracy more than ABC and vice versa. Training with both leads to a superior performance on all evaluated datasets. The data from the last iterations, UT_6 and ABC_3 , further improve the accuracy with an average improvement of 3.7 % relative to OCR_b .

The model trained with the UT'_1 and ABC'_1 datasets is also evaluated. Those datasets are subsets of the UT_1 and ABC_1 datasets that would have been obtained if no neighbourhood search or edit distance filtering was used. With the exception of IC03 dataset, this model’s performance is always inferior to the model trained with all the data.

The results also show that while the baseline model, trained on synthetic data only, performs well over a wide range of datasets, the performance on UT is rather poor - only 52.8 % accuracy. This shows the challenging nature of the dataset, partially due to the presence of heavily blurred images and the high frequency of vertical/diagonal text direction. The summary of the experiments can be seen in Table I.

For comparison with other methods, we trained and evaluated our best performing model on alpha-numeric characters only. The baseline model is pretrained with MJ and ST and fine-tuned with the the UT_6 and ABC_3 . During evaluation, all images with unknown characters are filtered out. The boost in performance here is slightly lower, 3.3 % on average.

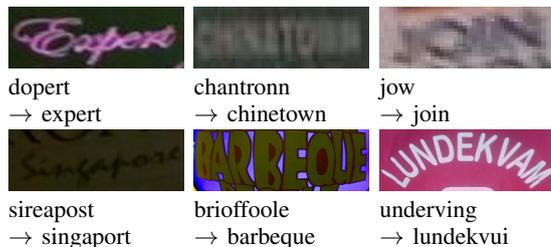


Fig. 7: Images with improved results after PGT training. The original OCR_b and $OCR_{UT_6+ABC_3}$ (\rightarrow) predictions.



Fig. 8: Images with worse results after PGT training. The original OCR_b and $OCR_{UT_6+ABC_3}$ (\rightarrow) predictions.

This single model achieves second-best performance on three different datasets with respect to the most recent state-of-the-art models. Note that state-of-the-art performance is achieved by different architectures trained on different data across the benchmarks, furthermore, the model of [41] uses multiple different model configurations.

The results of our work show that training with automatically generated PGT from very different domains, such as born-digital documents, can significantly improve the performance of a recognition model over a wide range of scene-text datasets. Also, adding only a relatively small number of those images helps significantly, implying the variety of data is important. Some common characteristics of the images where the PGT data has improved the model’s performance are blurred images, perspective distortions, artistic/handwritten fonts or occluded/cropped characters. Examples are shown in Figure 7 while images where the performance has deteriorated are shown in Figure 8.

VII. CONCLUSION

We proposed a PGT generation method and applied it to two different sources of weak annotations. As our baseline model, we chose the best-performing architecture with publicly available code-base [36]. Training with PGT without any architecture changes consistently improved the recognition accuracy both on images from the same domain and across different benchmark datasets and thus different domains. Our method is architecture-agnostic and thus can be applied to improve the performance of other architectures as well.

REFERENCES

- [1] Veit *et al.*, “Coco-text: Dataset and benchmark for text detection and recognition in natural images,” *arXiv:1601.07140*, 2016.
- [2] Chen *et al.*, “Cross-domain scene text detection via pixel and image-level adaptation,” in *NIPS*. Springer, 2019, pp. 135–143.

	Training dataset - % in batch					Evaluation on							Summary			UT _w	Δ
	Full			Weak		IIIT	SVT	IC03	IC13	IC15	SP	CT	Δ				
	MJ	ST	MLT	UT	ABC								avg	min	max		
OCR _b	45	45	10	0	0	89.8	86.7	94.3	91.2	68.5	77.2	72.3	0.0	0.0	0.0	52.8	0.0
OCR _{ABC₁}	30	30	10	0	30	92.8	88.9	94.8	93.0	71.0	77.5	73.6	1.7	0.3	3.0	53.9	1.1
OCR _{UT₁}	30	30	10	30	0	92.2	88.6	95.0	94.1	70.6	79.2	76.4	2.3	0.7	4.1	61.6	8.8
OCR _{UT₁+ABC₁}	20	20	10	20	30	93.0	89.2	95.2	94.1	71.6	79.2	77.8	2.9	0.9	5.5	61.8	9.0
OCR _{UT₆+ABC₃}	20	20	10	20	30	93.5	90.7	95.5	94.0	74.6	80.1	77.8	3.7	1.2	6.1	67.8	15.0
OCR _{UT₁+ABC₁}	20	20	10	20	30	91.4	88.1	95.5	93.9	69.5	77.1	74.0	1.4	-0.1	2.7	59.1	6.3

TABLE I: Recognition rate on standard benchmarks, non-alphanumeric characters included. Validation was performed on the union of training sets, with the exception of the UT dataset. Average, min. and max. improvements relative to OCR_b (Δ). UT_w results are not added to the summary since PGT was extracted from the UT domain and higher improvements are observed.

	Training dataset - % in batch					Evaluation on							Summary			UT _w	Δ
	Full			Weak		IIIT	SVT	IC03	IC13	IC15	SP	CT	Δ				
	MJ	ST	UT ₆	ABC ₃	avg								min	max			
OCR _b	50	50	0	0	87.6	88.6	94.5	91.9	75.2	78.9	73.2	0	0	0	54.9	0.0	
OCR _{UT₆+ABC₃}	25	25	20	30	91.7	91.8	95.7	94.2	80.0	82.5	77.4	3.3	1.2	4.8	68.9	14.0	
Published SOTA References	SOTA methods were each trained on different datasets					95.3	92.7	96.6	96.4	82.8	87.0	88.5					
SOTA 2 nd References						[15]	[41]	[41]	[42], [15]	[41]	[41]	[15]					
						94.4	91.8	95.4	94.7	78.7	83.6	87.5					
						[43]	[15]	[44]	[41], [44]	[43]	[43]	[41], [43]					

TABLE II: Recognition rate on standard benchmarks, non-alphanumeric characters excluded. Like many published methods, the model was pre-trained with MJ and ST datasets. Average, min. and max. increment relative to OCR_b (Δ). The second best scores and our score within 1 %, 3 %, 5 % from SOTA highlighted.

- [3] R. Gomez *et al.*, “Selective style transfer for text,” in *ICDAR*, 2019, pp. 805–812.
- [4] Zhan *et al.*, “Ga-dan: Geometry-aware domain adaptation network for scene text detection and recognition,” in *ICCV*, 2019, pp. 9105–9115.
- [5] Gupta *et al.*, “Synthetic data for text localisation in natural images,” in *CVPR*, 2016, pp. 2315–2324.
- [6] Liao *et al.*, “Synthtext3d: Synthesizing scene text images from 3d virtual worlds,” *arXiv:1907.06007*, 2019.
- [7] S. Long and C. Yao, “Unrealtext: Synthesizing realistic scene text images from the unreal world,” *arXiv:2003.10608*, 2020.
- [8] Baek *et al.*, “Character region awareness for text detection,” in *CVPR*, June 2019.
- [9] Sun *et al.*, “Chinese street view text: Large-scale chinese text reading with partially supervised learning,” in *ICCV*, October 2019.
- [10] Qin *et al.*, “Curved text detection in natural scene images with semi-and weakly-supervised learning,” *arXiv:1908.09990*, 2019.
- [11] Zhang *et al.*, “Uber-text: A large-scale dataset for optical character recognition from street-level imagery,” in *Scene Understanding Workshop-CVPR*, 2017.
- [12] Epshtein *et al.*, “Detecting text in natural scenes with stroke width transform,” in *CVPR*. IEEE, 2010, pp. 2963–2970.
- [13] L. Neumann and J. Matas, “A method for text localization and recognition in real-world images,” in *ACCV*. Springer, 2010, pp. 770–783.
- [14] Jaderberg *et al.*, “Reading text in the wild with convolutional neural networks,” *IJCV*, vol. 116, no. 1, pp. 1–20, 2016.
- [15] Liao *et al.*, “Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes,” *IEEE t. PAMI*, 2019.
- [16] Long *et al.*, “Textsnake: A flexible representation for detecting text of arbitrary shapes,” in *ECCV*, 2018, pp. 20–36.
- [17] Jaderberg *et al.*, “Spatial transformer networks,” in *NIPS*, 2015, pp. 2017–2025.
- [18] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv:1409.1556*, 2014.
- [19] He *et al.*, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [20] Graves *et al.*, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *ICML*, 2006.
- [21] Cheng *et al.*, “Focusing attention: Towards accurate text recognition in natural images,” in *ICCV*, 2017, pp. 5076–5084.
- [22] Liu *et al.*, “Fots: Fast oriented text spotting with a unified network,” in *CVPR*, June 2018.
- [23] Buřta *et al.*, “E2e-mlt-an unconstrained end-to-end method for multi-language scene text,” in *ACCV*. Springer, 2018, pp. 127–143.
- [24] Qin *et al.*, “Towards unconstrained end-to-end text spotting,” in *ICCV*, October 2019.
- [25] D.-H. Lee, “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks,” in *Workshop on challenges in representation learning, ICML*, vol. 3, 2013, p. 2.
- [26] Sun *et al.*, “Chinese street view text: Large-scale chinese text reading with partially supervised learning,” in *ICCV*, 2019, pp. 9086–9095.
- [27] Iwana *et al.*, “Judging a book by its cover,” *arXiv:1610.09204*, 2016.
- [28] Mishra *et al.*, “Scene text recognition using higher order language priors,” in *BMVC*, 2012.
- [29] Wang *et al.*, “End-to-end scene text recognition,” in *ICCV*. IEEE, 2011, pp. 1457–1464.
- [30] Trung *et al.*, “Recognizing text with perspective distortion in natural scenes,” in *ICCV*, December 2013.
- [31] Lucas *et al.*, “Icdar 2003 robust reading competitions,” in *ICDAR*. Citeseer, 2003, pp. 682–687.
- [32] Karatzas *et al.*, “Icdar 2013 robust reading competition,” in *ICDAR*. IEEE, 2013, pp. 1484–1493.
- [33] Karatzas *et al.*, “Icdar 2015 competition on robust reading,” in *ICDAR*. IEEE, 2015, pp. 1156–1160.
- [34] C. K. Chng *et al.*, “Total-text: Towards orientation robustness in scene text detection,” *IJDAR*, vol. 23, pp. 31–52, 2020.
- [35] Risnumawan *et al.*, “A robust arbitrary text detection system for natural scene images,” *Expert Systems with Applications*, vol. 41, no. 18, pp. 8027–8048, 2014.
- [36] Baek *et al.*, “What is wrong with scene text recognition model comparisons? dataset and model analysis,” *arXiv:1904.01906*, 2019.
- [37] Liu *et al.*, “Star-net: A spatial attention residue network for scene text recognition,” in *BMVC*, vol. 2, 2016, p. 7.
- [38] Shi *et al.*, “Robust scene text recognition with automatic rectification,” in *CVPR*, 2016, pp. 4168–4176.
- [39] Graves *et al.*, “Bidirectional lstm networks for improved phoneme classification and recognition,” in *ICANN*, 2005, pp. 799–804.
- [40] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [41] Litman *et al.*, “Scatter: Selective context attentional scene text recognizer,” *arXiv:2003.11288*, 2020.
- [42] Lu *et al.*, “Master: Multi-aspect non-local network for scene text recognition,” *arXiv:1910.02562*, 2019.
- [43] Yang *et al.*, “Symmetry-constrained rectification network for scene text recognition,” in *ICCV*, 2019, pp. 9147–9156.
- [44] Shi *et al.*, “Aster: An attentional scene text recognizer with flexible rectification,” *TPAMI*, vol. 41, no. 9, pp. 2035–2048, 2018.

APPENDIX A
BOUNDING BOX TRANSFORMATIONS

We define the following transformations to generate a set of bounding boxes in the neighbourhood of an input bounding box: Extending/shrinking the bounding box on the left/right/top. Angle modification and bottom extension/shrinkage were also considered but the benefits were insignificant. To keep the computational cost reasonable, we also assume that changes to the left side of the bounding box do not influence the recognition of the characters on the right side and vice versa, the optimization on each side is done independently. Horizontally, we extend/shrink the box with width w and height h in each direction by up to c characters, the character length being estimated as the average character length $ch_{avg} = \frac{w}{|text|}$. On the top, we extend by up to $\frac{h}{\beta}$ and shrink by up to $\frac{h}{\gamma}$.

Each of the transformed bounding boxes can be characterized by three integer parameters relative to the original bounding box - (t, l, r) - defining the extension/shrink (distinguished by the sign) by t, l, r units on top/left/right, where the horizontal unit is $\frac{ch_{avg}}{\delta}$ and the vertical unit is $\frac{h}{\kappa}$. Consequently, $l, r \in [-c \cdot \delta; c \cdot \delta]$ and $t \in [-\frac{\kappa}{\gamma}; 2\frac{\kappa}{\beta}]$. The transformed bounding boxes that exceed the image or do not overlap with the original one are immediately discarded.

To obtain the final bounding box bb_f characterized by (t_f, l_f, r_f) , we find an optimal bounding box in both directions (left and right). For each direction, we find the set of boxes $B = \{(t_i, l_i, r_i)\}_{i=1\dots n}$ that result in the lowest edit distance from g . The sets $T = \bigcup_{(t_i, l_i, r_i) \in B} t_i$, $L = \bigcup_{(t_i, l_i, r_i) \in B} l_i$ and $R = \bigcup_{(t_i, l_i, r_i) \in B} r_i$ are created.

From the boxes transformed in the left direction, we obtain

$$t_l = \min T \quad (10)$$

and

$$l_f = \frac{\min L + \min(\max L, o + \min L)}{2} \quad (11)$$

Analogously, for the right direction, we obtain

$$t_r = \min T \quad (12)$$

and

$$r_f = \frac{\min R + \min(\max R, o + \min R)}{2} \quad (13)$$

We set

$$t_f = \max(t_r, t_l) \quad (14)$$

In our experiments, the constants were assigned as follows: $c = 7$, $\beta = 2$, $\gamma = 4$, $\delta = 4$, $\kappa = 4$, $o = 8$. All of these but o control how many bounding boxes will be generated and the trade-off between precision of the box and computation time. They were selected ensuring all the generated boxes for one PGT proposal can be run in a single batch on GPU and by observing qualitative results. The $o = 8$ means that tt_f won't be more than one average character length wider on the right/left than the smallest bounding box that gives the correct recognition output.

	Iteration / Mined:	with neigh. s.	w/o neigh. s.	Δ
UT	1	92,909	72,990	19,919
	2	105,126	86,295	18,831
	3	109,557	90,480	19,077
	4	111,663	92,660	19,003
	5	113,046	93,994	19,052
	6	113,810	94,890	18,920
ABC	1	1,536,583	1,234,219	302,364
	2	1,581,109	1,354,219	226,890
	3	1,594,333	1,375,571	218,762

TABLE III: PGT generation on the Uber-Text training dataset where text location information is ignored. The number of captions generated in iterations 1 to 6, with and without the neighbourhood search, and the difference Δ .

	Training dataset - % in batch					Acc.	NED
	Full				Weak		
	MJ	ST	MLT	UT _F	UT (30)		
[11]						56.4	
OCR _b	45	45	10	0	-	41.6	35.2
OCR _{UT₁'}	30	30	10	0	UT ₁ '	55.2	28.0
OCR _{UT₁}	30	30	10	0	UT ₁	57.9	26.2
OCR _{UT₂}	30	30	10	0	UT ₂	62.7	23.3
OCR _{UT₃}	30	30	10	0	UT ₃	64.4	22.5
OCR _{UT₄}	30	30	10	0	UT ₄	65.4	21.5
OCR _{UT₅}	30	30	10	0	UT ₅	66.0	21.1
OCR _{UT₆}	30	30	10	0	UT ₆	66.1	21.2
OCR _{UT_F}	30	30	10	30	-	78.0	10.0
OCR _{UT_{PL99}}	30	30	10	30	UT _{PL99}	44.5	33.5
OCR _{UT_{PL90}}	30	30	10	30	UT _{PL90}	45.7	31.7
OCR _{UT_{PL80}}	30	30	10	30	UT _{PL80}	44.8	32.1
OCR _{UT_{PL50}}	30	30	10	30	UT _{PL50}	44.4	33.1

TABLE IV: Recognition rates and normalized edit distance (acc., NED) on the Uber-Text test set. The data obtained from the i^{th} iteration of the PGT generation is denoted as UT _{i} . UT_F and UT_{PL t} are the fully annotated and pseudo-labelled (with a threshold t) datasets, respectively. UT₁' is a subset of UT₁ obtained without the neighbourhood search.

APPENDIX B
ADDITIONAL EXPERIMENTAL RESULTS

In this section, we provide additional experimental results in detail.

The number of text instances located in each iterations of the PGT generation on both the UT and ABC datasets is reported in Table III.

The results of different models, tested and validated on the UT dataset, are reported in Table IV. For each model, we provide the per-batch percentage of examples from each dataset during training, the accuracy and the normalized edit distance.

A. *Semi-supervised learning via pseudo-labelling*

For comparison with prior work [10], [24], we implement the pseudo-labelling [25] approach to semi-supervised learning with confidence-thresholding of the pseudo-labels. Pre-trained models, in our case, OCR_b and TextSnake, are used to generate predictions on a set of unlabelled data - the UT training set (ignoring the existing ground truth). Only the predictions with a recognition confidence above a thresholds t are kept to retrain the recognition model.

Confidence threshold	Mined
0.99	100,242
0.90	137,005
0.80	154,560
0.50	204,549

TABLE V: The number of PGt data generated in the semi-supervised setup with different confidence thresholds. There were 365,200 detections in total.

We implement the method with different thresholds, $t \in \{99, 90, 80, 50\}$, the results are reported in Table IV. The best performing model is also included in Figure 6. The accuracy of this model is 12.2 % lower than the accuracy in the first iteration using our proposed method, OCR_{UT_1} .

APPENDIX C PGT ACCURACY

Different kind of errors and ambiguities can occur in the resulting PGT data. We report their numbers in Table VII, computed on 500 sample crops from the UT dataset and 500 samples from the ABC dataset.

The wrong text and not text errors are the most harmful ones. However, employing recognition confidence based filtering, these can be reduced significantly. This is because these errors mostly occur for very blurred texts or images with no text, where the network predicts a short, common word such as ‘the’, ‘in’, ‘on’, ‘at’ with low confidence.

The error classification is not always clear, some may overlap.

	Training dataset - % in batch					Evaluation on								Summary			UT _w	Δ
	Full			Weak		IIIT	SVT	IC03	IC13	IC15	SP	CT	Δ					
	MJ	ST	MLT	UT	ABC	3000	647	867	1015	2077	645	288	avg	min	max	20 000		
OCR _b	45	45	10	0	0	89.8	86.7	94.3	91.2	68.5	77.2	72.3	0.0	0.0	0.0	52.8	0.0	
OCR _{UT₁}	30	30	10	30	0	91.6	88.7	94.7	94.0	71.3	79.4	74.7	2.1	0.4	2.8	62.5	9.7	
OCR _{UT_{PL99}}	30	30	10	30	0	91.0	87.5	94.1	93.6	68.7	76.4	71.2	0.4	-1.1	2.4	55.9	3.1	
OCR _{UT_{PL90}}	30	30	10	30	0	90.6	87.2	94.0	92.5	68.5	77.8	71.2	0.3	-1.1	1.3	57.5	4.7	
OCR _{UT_{PL80}}	30	30	10	30	0	90.4	87.0	94.2	93.2	68.5	77.2	73.3	0.5	-0.1	2.0	56.7	3.9	
OCR _{UT_{PL50}}	30	30	10	30	0	90.9	88.4	94.3	93.4	69.1	76.9	72.9	0.8	-0.3	2.2	56.1	3.3	

TABLE VI: Results of OCR_{UT₁} are slightly different than the previous table because the models here were validated on uber-text.

Counted on 500 sample crops:	ABC	UT	PGT	Crop
wrong text	19	6	cancer	
not text	1	2	the	
ambiguous GT	14	14	java,	
unclear GT	5	20	the	
wrong weak label - punctuation	1	2	1:15,000	
wrong weak label - other	0	7	hilling services	

TABLE VII: PGT errors and ambiguities in the ABC₃ and UT₆ datasets. Counted on 500 sample crops from each.