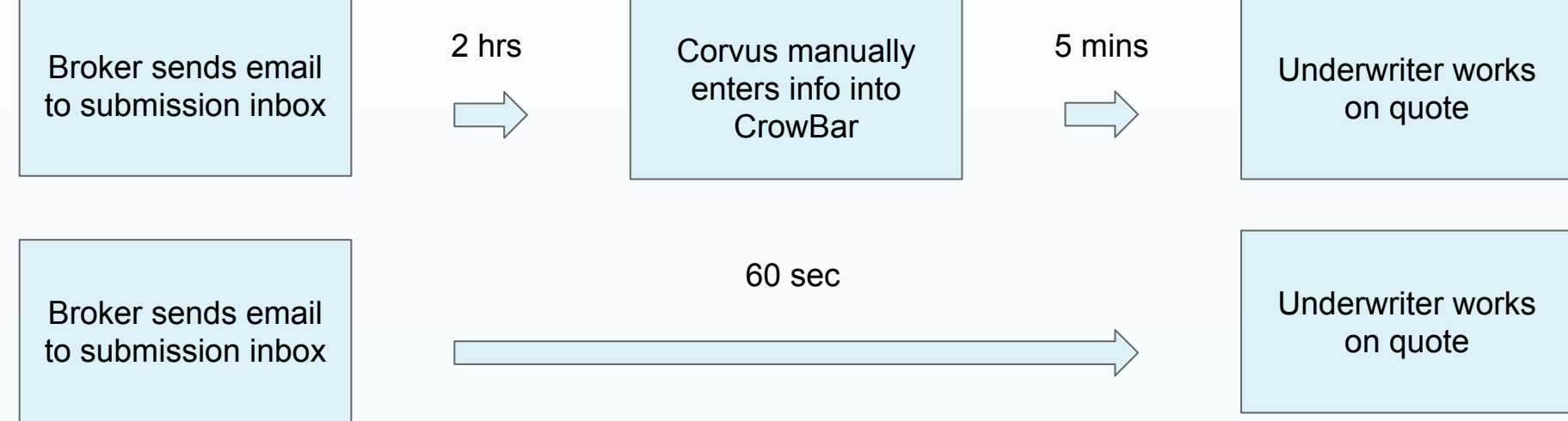


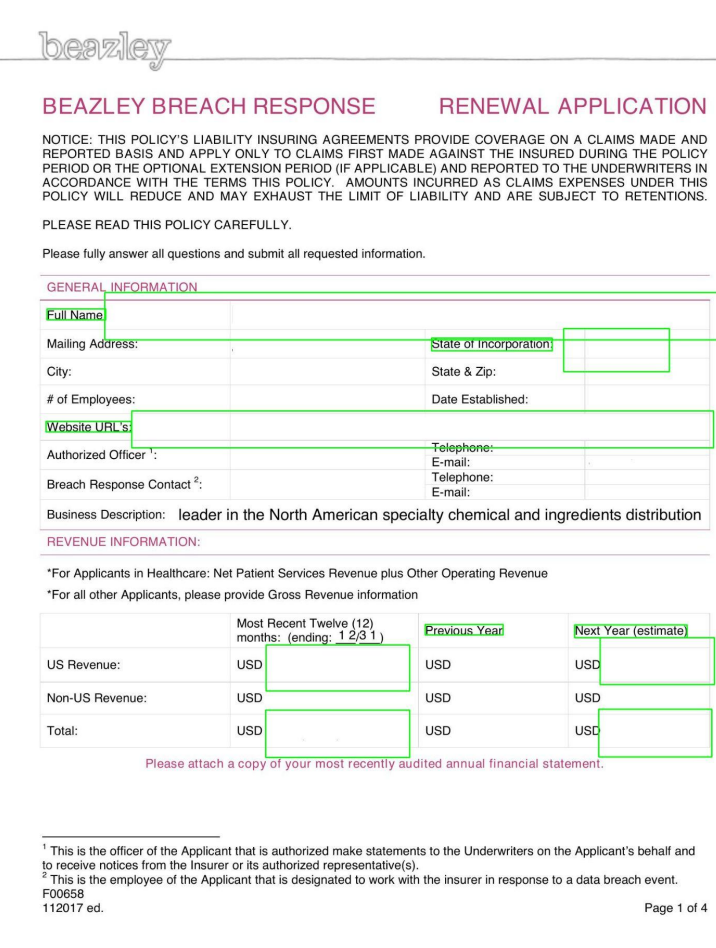
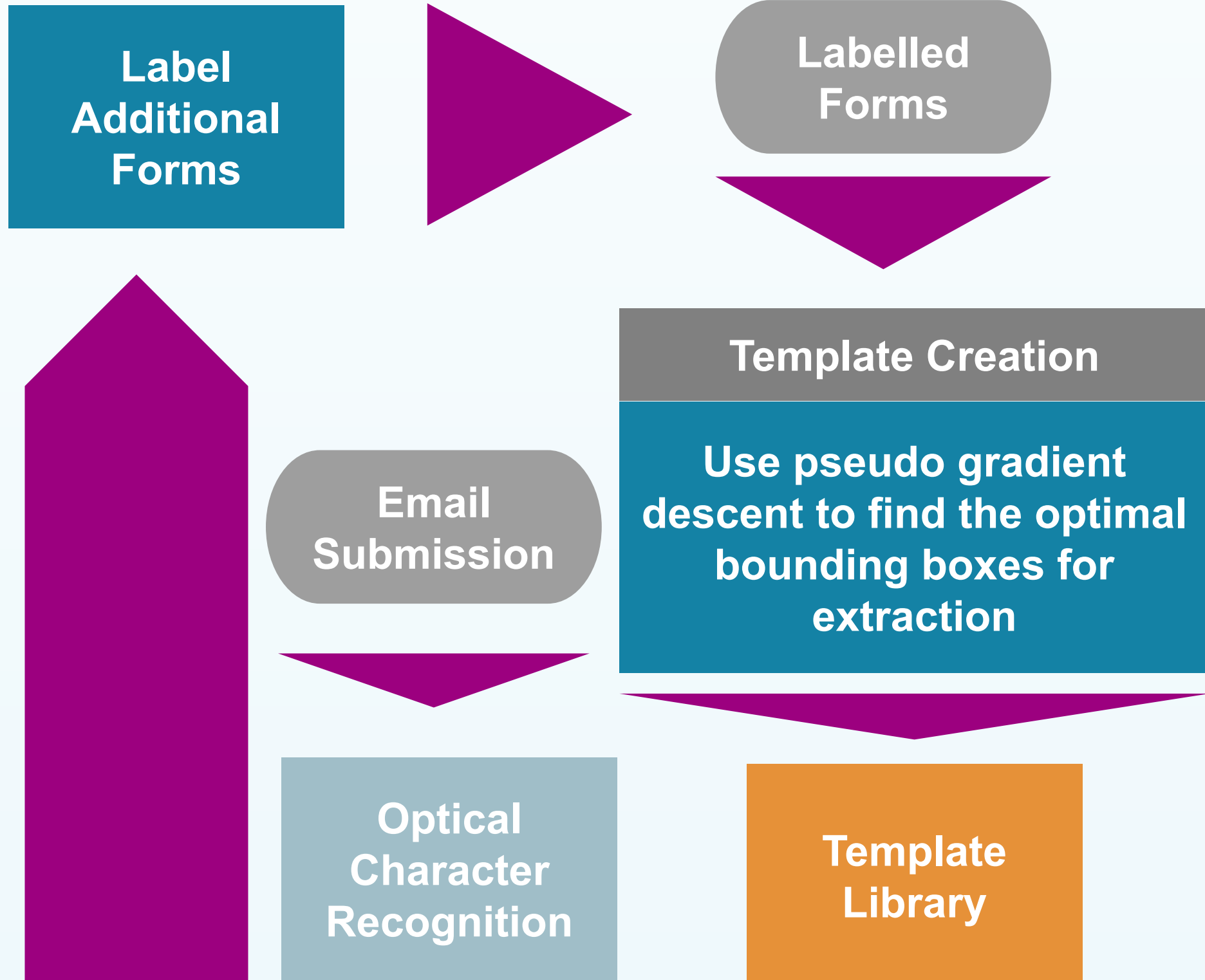
## Introduction

Corvus Insurance gets information for quote requests in two ways. The first way is brokers can directly enter the information into Corvus's proprietary software platform, the CrowBar. The second way is brokers send Corvus an email with the information, and this information typically comes in the form of PDF attachments. Currently, the process of taking the information from the PDF attachment and putting it into the CrowBar is done by an administrative support group. This intermediate step means that email quote requests will be slower than broker-entered ones. The goal of our project is remove the intermediate step by automating the data extraction of quote information from the PDFs.

**GOAL**



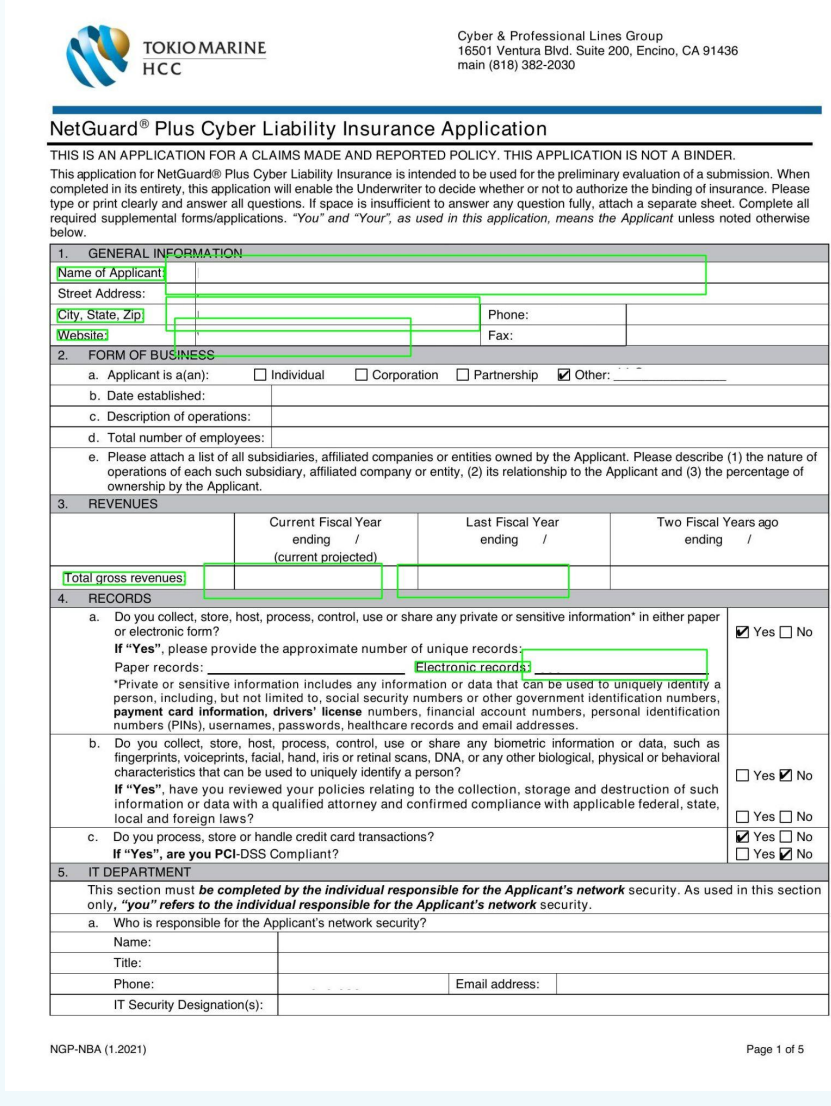
## Template Creation



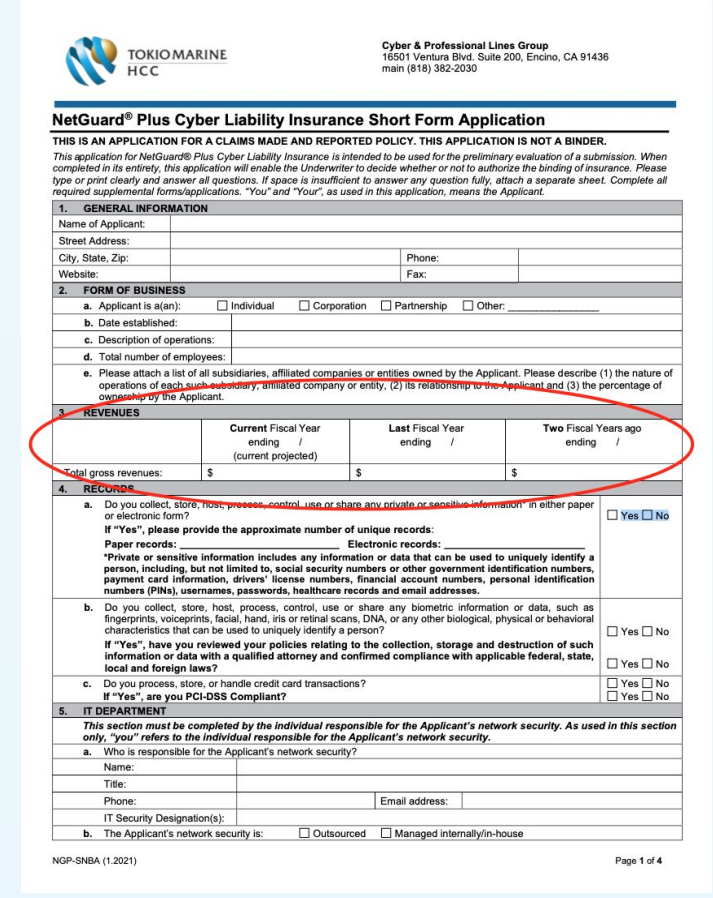
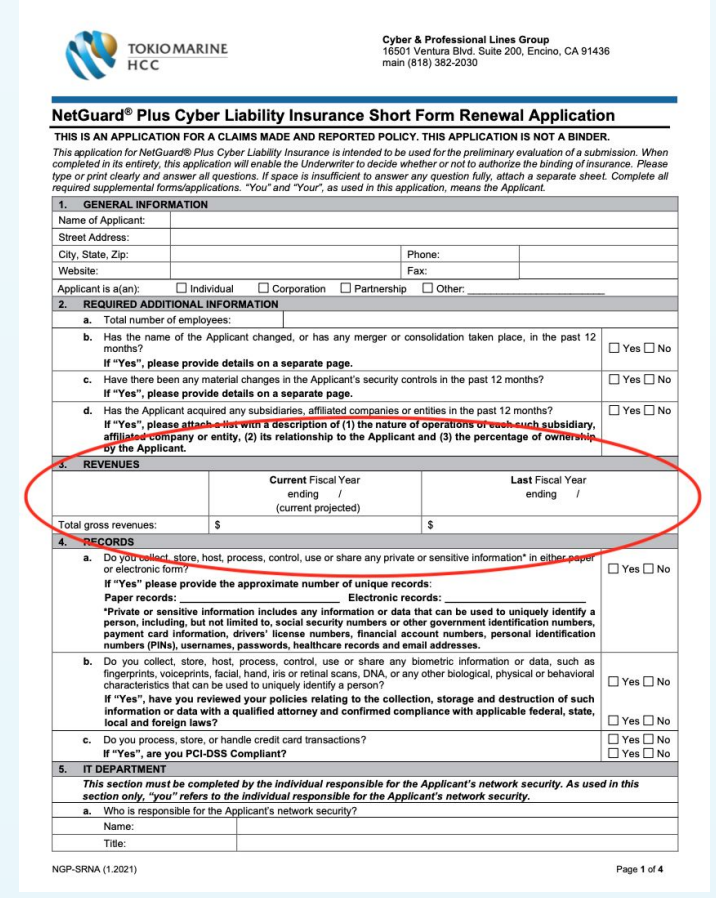
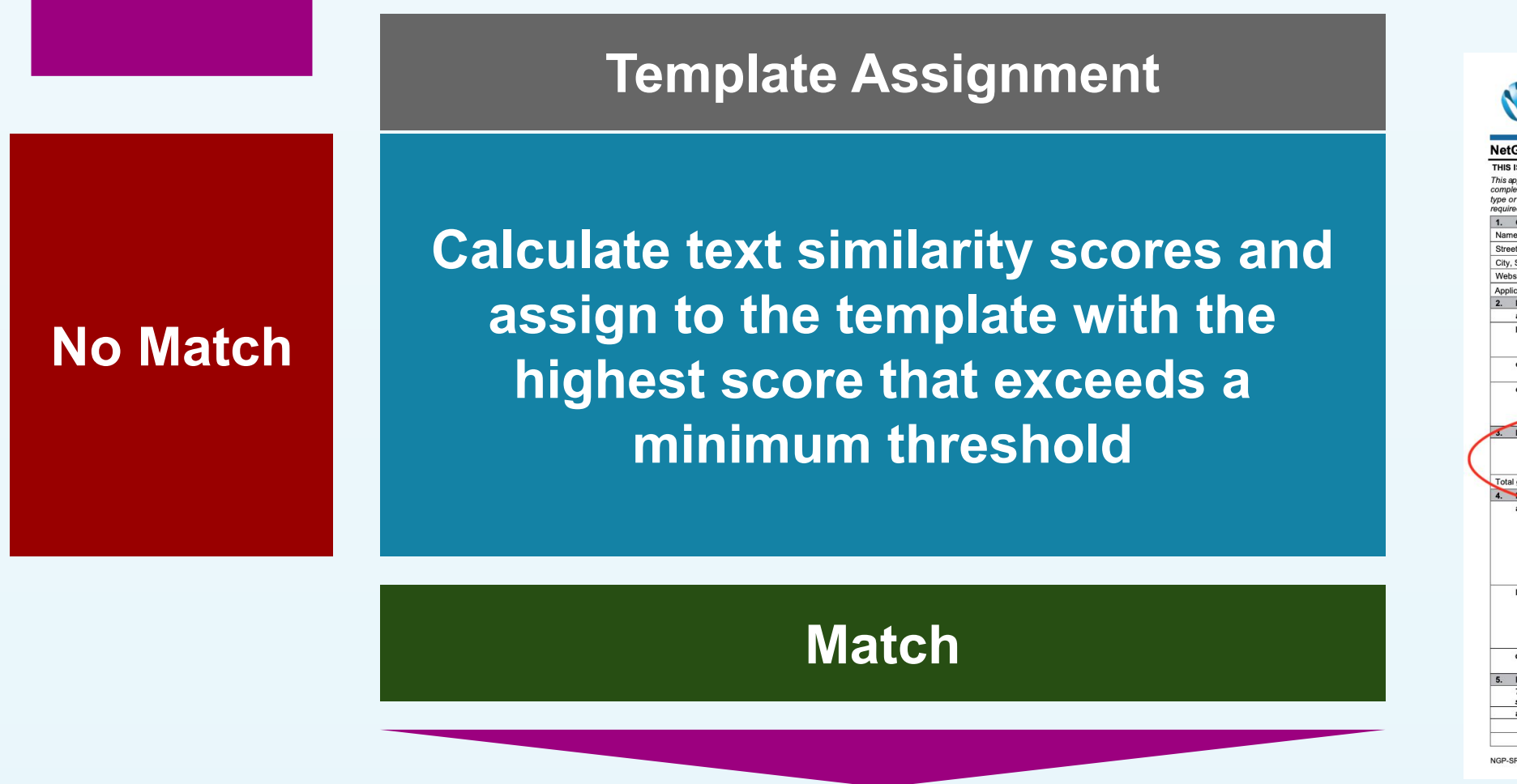
To extract data, we utilize templates that map the spatial relationship between a keyword phrase and the data to be extracted. We automated the template creation process. For each field we need to extract, we use a pseudo gradient descent algorithm to learn the spatial relationship between the keyword phrase and the associated value.

To do this, we need a set of similar forms that should belong to the same template. We did this by picking one form and querying our PDF database to find similar forms scored by cosine similarity. This is the same process used in the template identification that is elaborated on in the next section.

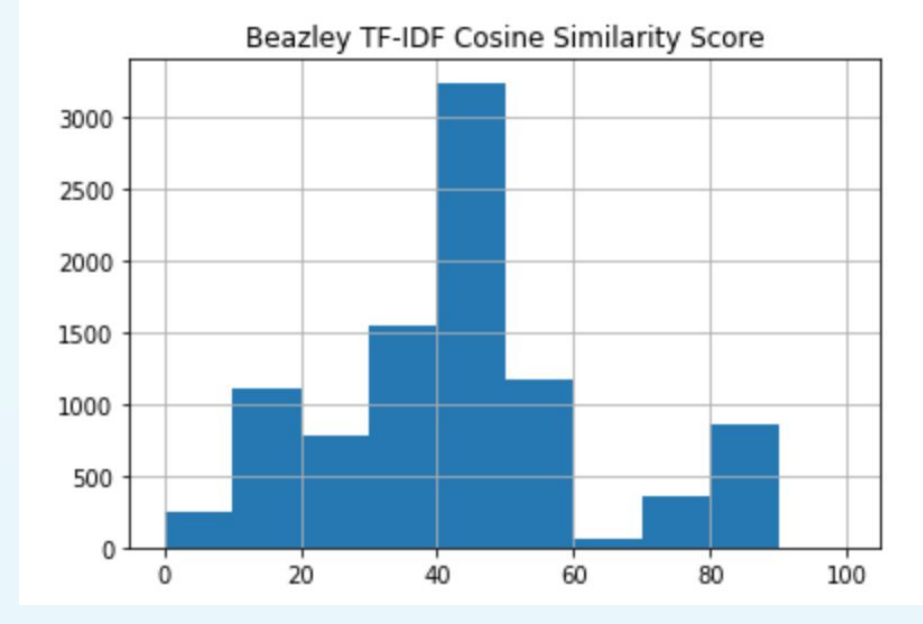
Optimization occurs by shifting the dimensions of the box, looking at the text within, and comparing it to manually labeled text using string-matching. Intuitively, the algorithm stacks each form on top of each other and selects the box that will capture the correct text for the most forms.



## Template Assignment



The assignment algorithm needs to differentiate between similar looking forms that have material differences in the location of the text to be extracted. For example, there is a key difference in the revenue section of the otherwise similar forms on the left.



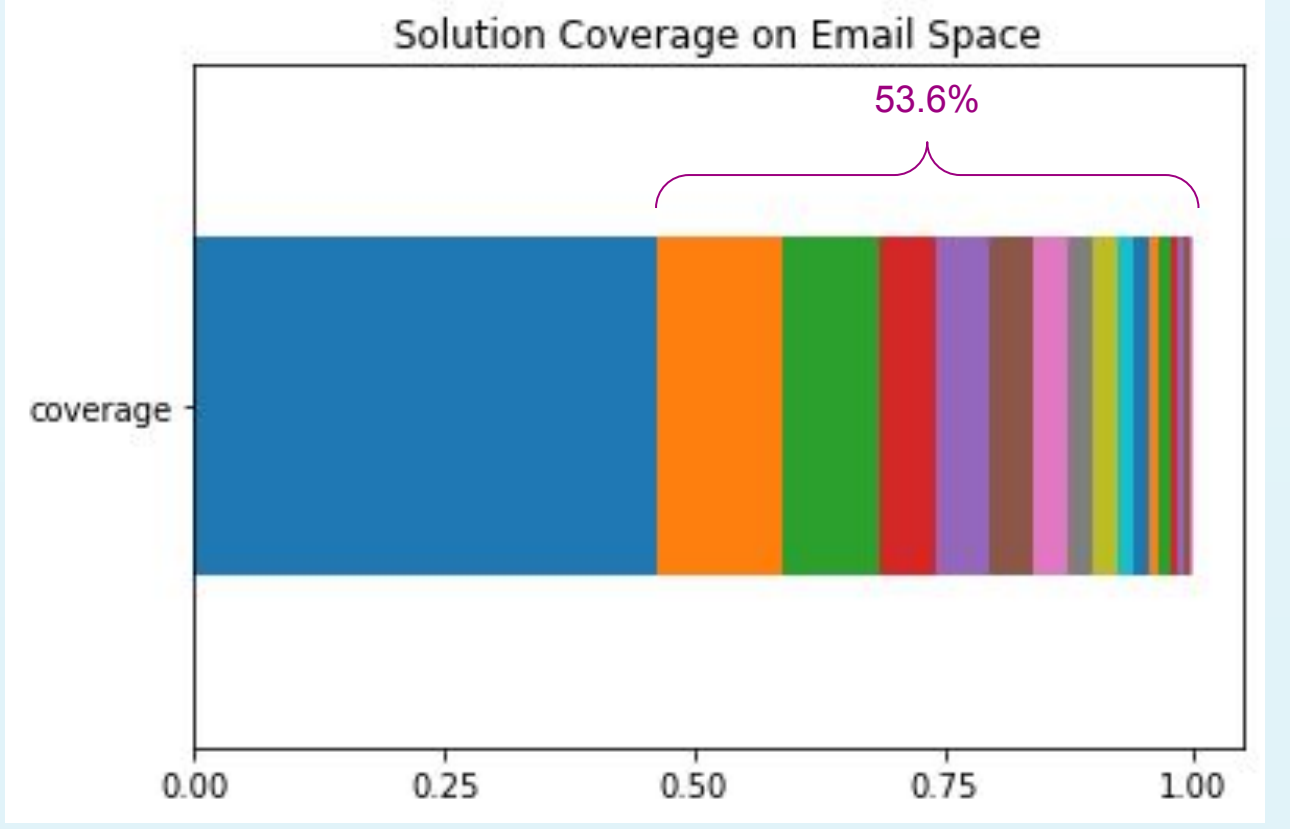
We created a vector representation of each form using term frequency - inverse document frequency (TF-IDF) for all of the text recognized by the optical character recognition engine on the first page of the form. We created a similarity score using cosine similarity scores of each vector representation.

Each form was assigned to the template with the highest similarity score so long as that score exceeded a minimum threshold. The Threshold for the Beazley template (with scores shown above) was 75.

## Results

The performance of the extraction process can be measured on two dimensions: the percentage of email submissions that are assigned to a template (coverage), and the accuracy of the predictions made by each template.

**Coverage:**  
We achieve 53.6% coverage on the cyber email submissions using the templates that we created thus far (about 20).



Here we show what each template contributes to entire space. All the space under the bracket represents the amount of emails that are covered by our solution.

**Accuracy:**  
We get correct predictions for all five target fields on over 70% of submissions for two of the most common templates.

Common sources of error are:  

- The optical character recognition engine interprets the text incorrectly. This is most common for handwritten forms and forms with poor scan quality
- The extraction bounding box is the incorrect size
- The wrong template is assigned.

Template	Forms Evaluated	All Correct	Accuracy
Beazley	51	74.5%	89.4%
Travelers	50	72%	90.8%
Tokio Marine	52	56%	84.7%
Combined	153	67.8%	88.3%

## Impact

**Scalability:**  
Corvus insurance is growing rapidly. A human administrative support team cannot scale at the rate Corvus expects their email submissions to grow. This extraction process by contrast will scale easily.

**Time to Quote:**  
The first company to reply to a quote request is an order of magnitude more likely to win the business. This solution will greatly reduce the time to quote.

**Potential to Create Novel Datasets:**  
The current application of this solution is to collect basic data for quote requests. However, the solution can be easily generalized to generate other data sets from PDF files. Corvus Insurance has expressed interest in using the solution for this purpose.

**Cost Savings:**  
The most immediate benefit to Corvus Insurance comes from time savings for the administrative support team currently performing data entry. Corvus estimates this will save tens of thousands of dollars. This number will increase dramatically as Corvus scales.

The second two errors may be reduced as additional forms are labelled.