# Evaluation Criteria for the *PatentsView Inventor Disambiguation Technical Workshop*

American Institutes for Research

July 22, 2015

Algorithms will be evaluated in two phases. The result of the first phase will be used to determine which groups will be invited to participate in the second phase. Due to resource constraints, we may not be able to invite all workshop participants to continue on to the second phase.

## 1  First Phase

This is an initial round of self-testing where participants will infer links for the bulk patents database. They may train their algorithms using any part of the provided data, as well as any additional data sets that have been submitted to the workshop organizers. However, participants should keep in mind that during the second phase of evaluation they will be asked to train and evaluate their algorithm on different subsets of labeled data.

The output file should be a tab-delimited file with two columns and no header. The first column should be an inventor ID that is constructed taking the hyphenated combination of the patent number and sequence fields for each inventor in the rawinventor table. The second column should be an integer ID generated by your program. Inventor IDs that are predicted to refer to the same unique individual should be assigned the same integer ID. For example, in the following excerpt the second author on the first patent and the first author on the second patent are believed to be the same individual:

```
1234567-1 1
1234567-2 2
1234567-3 3
2345678-1 2
2345678-2 4
```

In the first phase, algorithms will be evaluated on the following two criteria:

- Algorithm accuracy

- Run-time

For this end, elements will be taken into account to assess the performance of the algorithms: **recall rate**, **precision rate**, and **self-reported runtime of the algorithm**.

## 1.1 Recall Rate

Recall rate[1] is the proportion true links that were correctly predicted by the algorithm. The mathematical definition of recall rate is:

$$\text{Recall} = \frac{\text{\# of true positives}}{\text{\# of true positives} + \text{\# of false negatives}}$$

Note that the inventor disambiguation literature uses two statistics closely related to the recall rate to assess the performance of algorithms: splitting rate and lumping rate. Ventura et al. (2015) define the splitting rate as follows[2]:

$$\text{Splitting} = \frac{\text{\# of false negatives}}{\text{\# of true positives} + \text{\# of false negatives}}$$

The splitting rate is thus the complement of the Recall rate. That is:

$$\text{Recall} = 1 - \text{Splitting}$$

Ventura et al. (2015) define the lumping rate as follows:

$$\text{Lumping} = \frac{\text{\# of false positives}}{\text{\# of true positives} + \text{\# of false negatives}}$$

Although the relationship between the recall rate and the lumping rate is not as straightforward as it is between the splitting rate and the recall rate, a simple transformation gives the following equation:

$$\text{Recall} = \frac{\text{\# of true positives} + \text{\# of false positives}}{\text{\# of true positives} + \text{\# of false negatives}} - \text{Lumping}$$

## 1.2 Precision rate

The precision rate[3] is the proportion of true links from all predicted links. The mathematical definition of precision rate is:

$$\text{Precision} = \frac{\text{\# of true positives}}{\text{\# of true positives} + \text{\# of false positives}}$$

It is important for accuracy to understand what the algorithm?s performance is when trained on a dataset with different feature distributions or different feature importance in determining matches from those present in the test data sets. Because of this, we will compute precision and recall rates for the entire set of labeled data (including labeled data that was previously withheld from the released training data set), as well as for different biased samples. Judges will assess these statistics during this phase of evaluation.

---

[1]https://en.wikipedia.org/wiki/Precision_and_recall

[2]Please note that Ventura et al. 2015 versions of the splitting and lumping rates are different from those used by other researchers.

[3]https://en.wikipedia.org/wiki/Precision_and_recall

### 1.3 Run time of algorithm

Workshop participants will be asked to self-report the run time of their algorithms as well as the computing setup they used. Judges will take this information into consideration when evaluating the algorithms. Participants should take the following guidelines into account:

- The algorithm should not run for more than 5 days when processing all patent application and grant data (2001-2014 for applications; 1976-2014 for grants)

- The implementation should be runnable on hardware equivalent to a single Amazon Web Services (AWS) instance. For reference, currently the largest compute-optimized AWS instance provides 36 virtual CPUs and 60 GB memory.

- AIR and the panel will review any requests for software or hardware updates that might be required to accommodate the incorporation of a novel algorithm into the current PatentsView workflow. These requests must be submitted in your letter of intent to participate.

## 2 Second Phase

In the second phase, algorithms will be evaluated on the following three criteria:

- Algorithm generalization

- Run-time

- Usability of the implementation

The main goals of this stage are to reproduce the results of the first phase in a controlled environment, and to continue testing the generalizability of the algorithms. To do this, participants will be asked to re-run the algorithms on the entire patents database in a server environment that we will provide (participants will not be asked to inference the full patents database more than once during this phase). To test the generalizability of the algorithm, we will provide participants with new subsets of labeled data on which to train their algorithms, as well as non-overlapping sets of labeled data on which to test them. Recall and precision rates - previously described -will be computed and will be taken into account by the judges. The output format is the same as the output format in the first phase.

In this phase, all competitors should run their code on the provided server environment and should also provide the judges user documentation about how to run the program. This documentation will be used in part to evaluate the usability of the algorithm implementation. Usability will be taken into account during the final evaluation, alongside with the performance indicator and runtime.

Participants will self-report the algorithm runtime.