# Prescription database construction by data extraction from scanned files
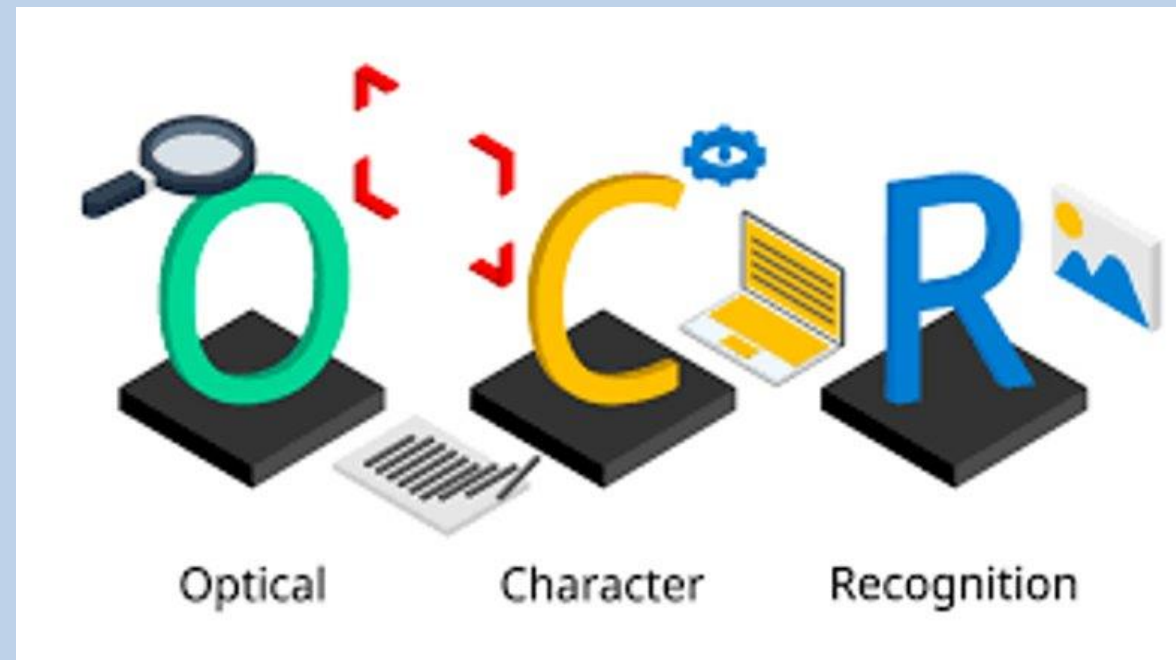
[1]Noah Hallberg, [1]Ryan Kwong, [2]Ankur Malik, [1]Saul Means, [1]Udayan Pandey, [1]Bharath Sadagopan, [2]Kabir Snell, [1]Varon Srinivasan, [1]Margaret Wang

[1]Purdue University, [2]University of California Santa Barbara

**PURDUE UNIVERSITY** | The Data Mine

**Inogen**

## Introduction

- Our corporate mentors tasked the joint Purdue and UCSB research group this semester to help them evaluate different methods in analyzing bulk data
- The goal was to help them chose the most optimal OCR solution for their needs
- Optical Character Recognition (OCR): using computer models to analyze data that tends to be handwritten
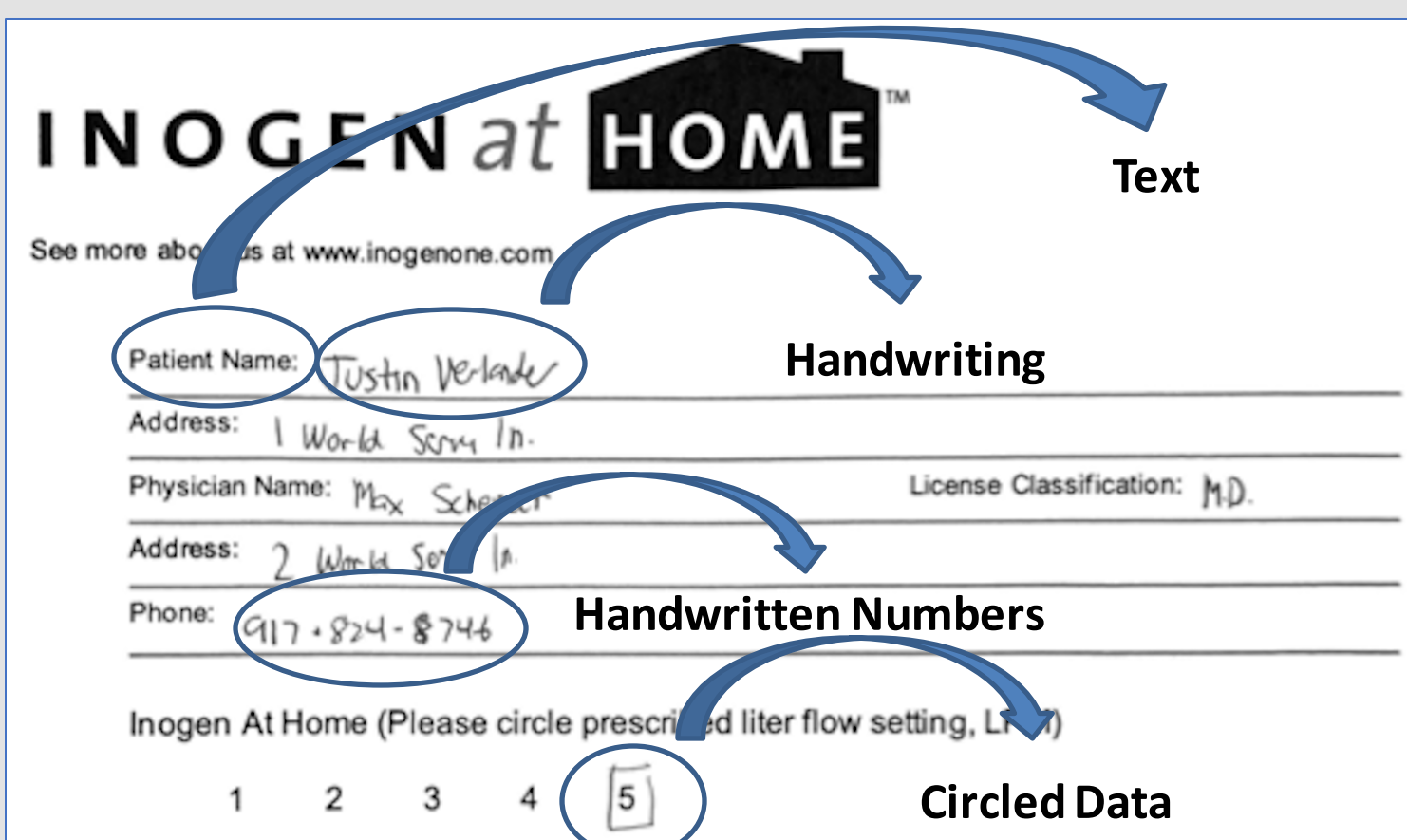


Optical   Character   Recognition

## Methodology

Overview
- Creating our own mock database
- Criteria for choosing a software
- Testing the software

Three main software considerations:

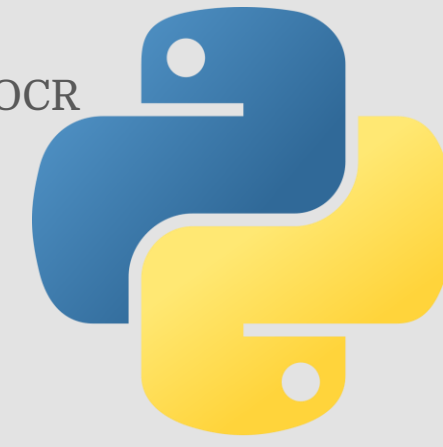1. HIPAA compliance
2. Performance
3. Price



**INOGEN at HOME**

See more about us at www.inogenone.com

Patient Name: Justin Verlade
Address: 1 World Serv In.
Physician Name: Max Scherzer       License Classification: M.D.
Address: 2 World Serv In
Phone: 917 · 824 · 8744

Inogen At Home (Please circle prescribed liter flow setting, L Min)

1   2   3   4   ⑤   6

Text
Handwriting
Handwritten Numbers
Circled Data

*No Inogen customer or patient data was used in this visual. Hypothetical patient information used for example purposes only.*

- Evaluated each on a scale of 0 to 1
- 0 = failed entirely
- 1 = perfect performance

## Evaluated Methods

Tesseract OCR (Python) :
- Package in python that implements Google's OCR software

Pros:
- Free (Open-Source Software)
- Highly customizable
- Very fast

Cons:
- Lower accuracy on OCR tests
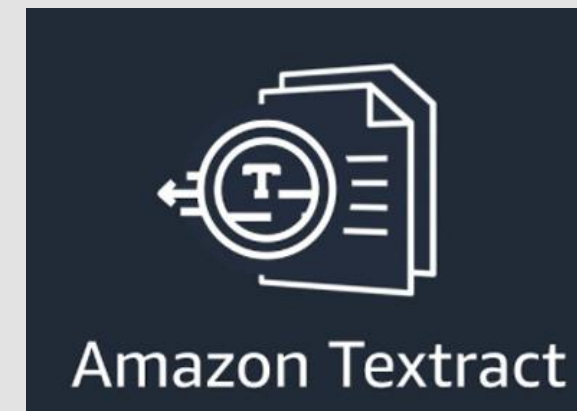- Collaboration becomes more complicated

Version 3.0

**UiPath™**

Automation platform with OCR Component (OCR scripting)
Pros:
- ~80% handwriting accuracy
- Effectively extracts to a CSV
Cons:
- Failure to read circles
- Requires a large amount of processing power

Version 2022.10.5

Machine Learning Optical Recognition Service (AWS)
Pros:
- Very good at handwriting
- Adaptable to numerous form types
Cons:
- Base version bad with circles
- Sometimes interprets scribbled out characters

**Amazon Textract**

2023 Version

**Base64.ai**
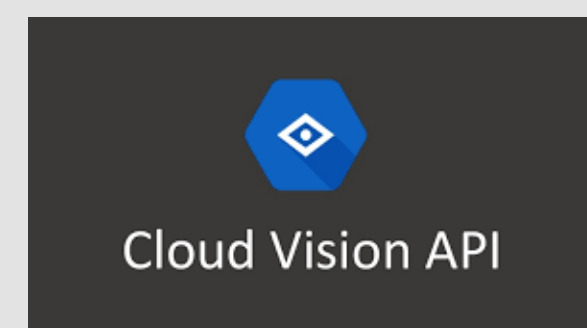
Artificial Intelligence OCR software
Pros:
- Handwriting and circle/box filling
- Easy, Support from Base64.
Cons:
- Circled Data
- Some issues with tables

Version 2023

OCR software
Pros:
- Business oriented model
- Trainable
Cons:
- Inferior performance compared to other software tested
- Does not limit outputs to a lexicon of English words or the Latin alphabet

**ABBYY®**

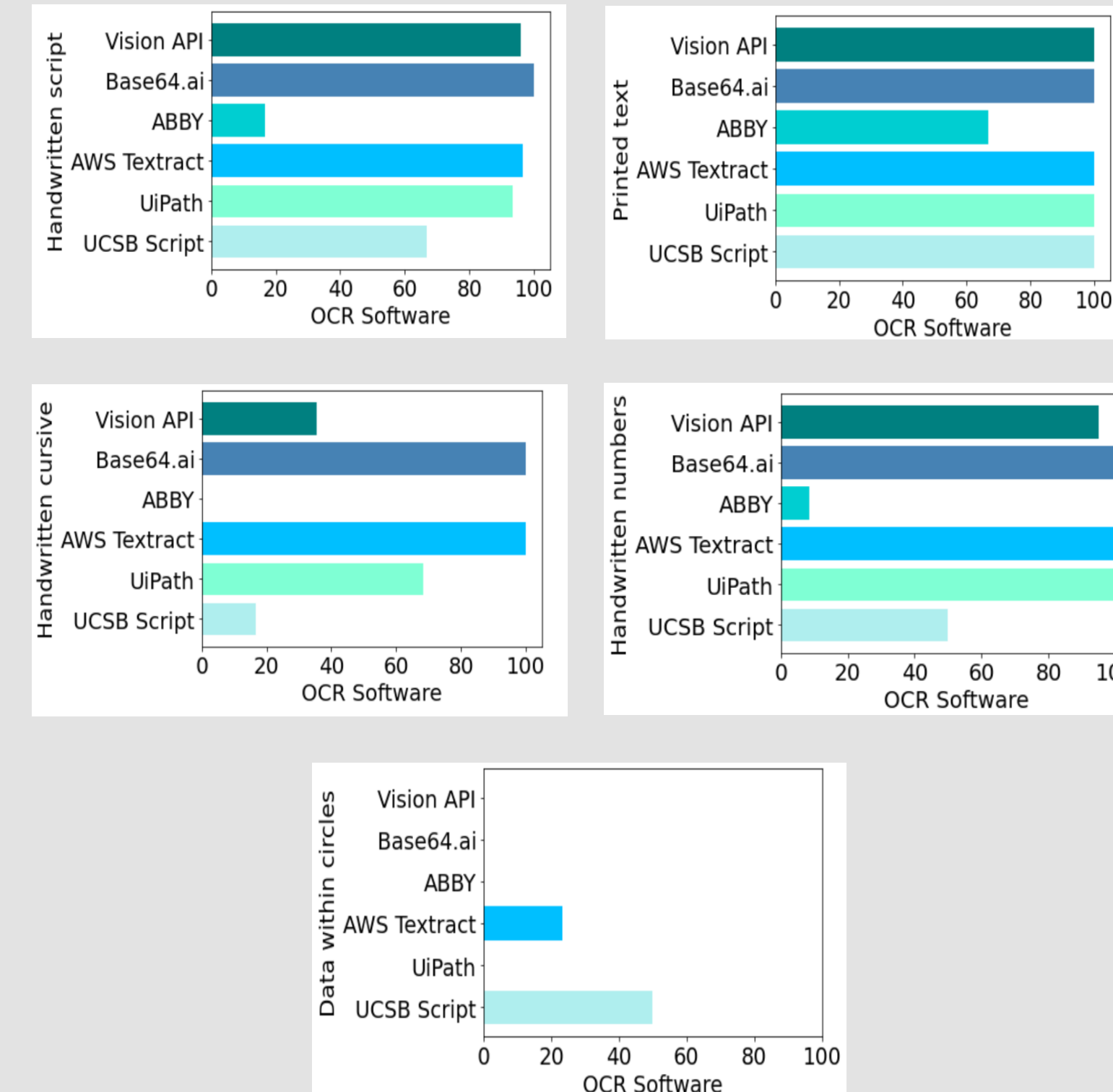Fine Reader 16
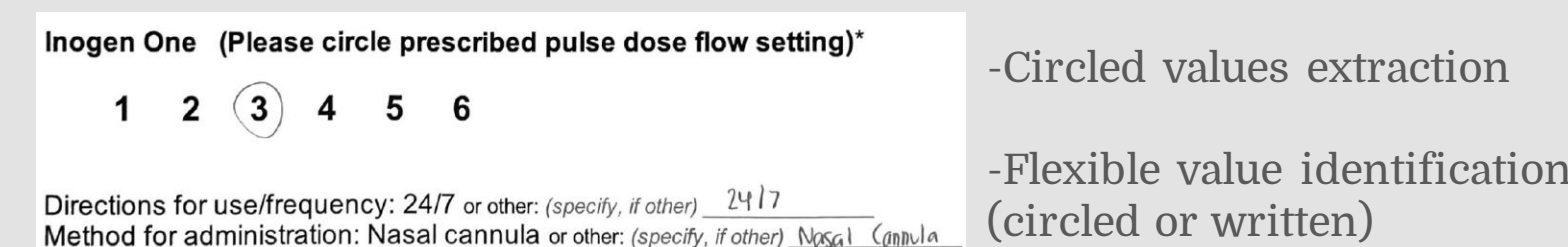
**Cloud Vision API**

Document AI – machine learning, OCR platform
Pros:
- Confident text extraction
- Identified symbols and room for ML
Cons:
- Poor results with circle data
- Unclear data sorting

Version 2023

## Combined Results (%Accuracy)



## Current Work

Working with companies to get customized solutions:



Inogen One (Please circle prescribed pulse dose flow setting)*

1   2   ③   4   5   6

Directions for use/frequency: 24/7 or other: (specify, if other) 24/7
Method for administration: Nasal cannula or other: (specify, if other) Nasal Cannula

-Circled values extraction

-Flexible value identification (circled or written)

Working on improving Python solution:

- Testing new packages
- Identifying regions of interest
- Improving hand-writing accuracy
- Image pre-processing

## References

https://docs.aws.amazon.com/textract/index.html
https://guides.nyu.edu/tesseract/usage
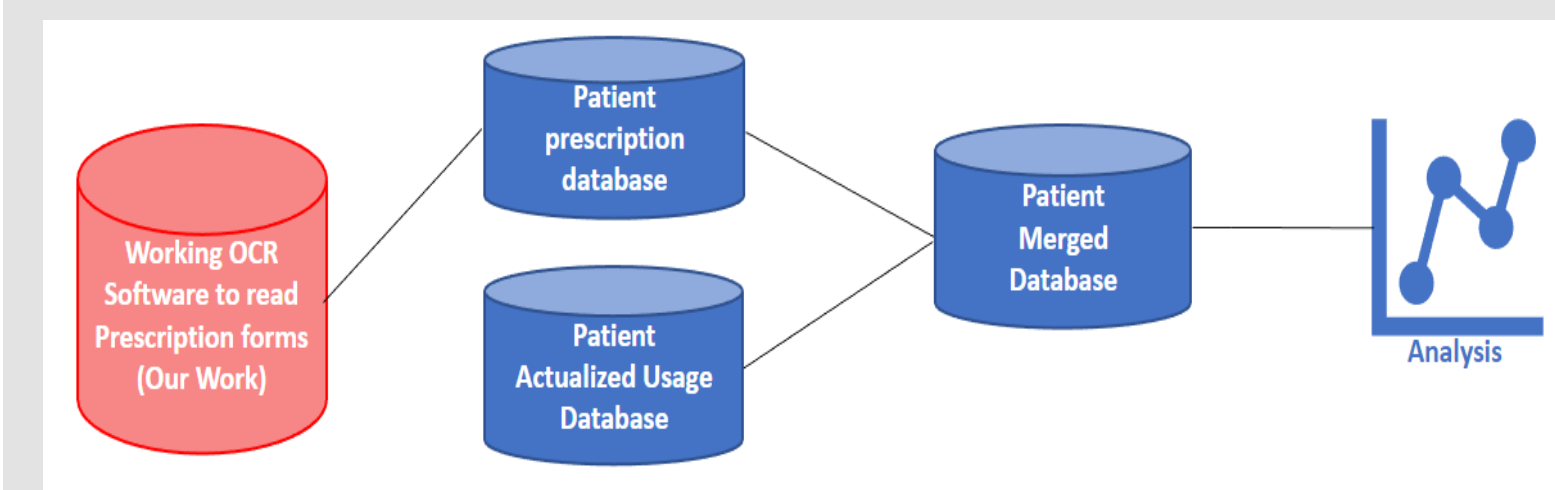https://base64.ai/features/data-extraction-api/handwritten-document

## Conclusions

- Created and tested mock patient prescription PDFs
- AWS and Base64.ai were the most accurate and flexible software
- Tesseract (python) and Google API can be explored
- Circled values weren't recognized by any method

## Future Goals

We are creating a tool that allows Inogen to generate a patient prescription database
This new database can help them answer questions like:
- What proportion of patients are adhering to their prescribed flow setting?
- What kind of patients are not adhering to their prescribed flow setting?
- In cases of inconsistency with flow setting and prescription; how is the flow setting being misused (higher or lower)?



## Acknowledgements