

AUGUST 2018

Handwriting Recognition POC

Google Cloud Platform API Assessment



Executive Summary

- A proof of concept is needed to determine the feasibility of **using Google's Cloud API services for text detection and optical character recognition (OCR) applied to common claims forms**
- Google Cloud Platform (GCP) offers a number of pre-trained and available machine learning (ML) models, including the Cloud Vision API, capable of out of the box image labeling, logo detection, OCR, and text detection.
- The proof of concept will aim to **understand the capabilities and limitations of the Cloud Vision API**, and suggest a potential ML pipeline for automation of claim intake.
- The success criteria for the proof of concept will be defined as follows:
 - Determine which available APIs and models within GCP are applicable to automation of claims intake
 - Using determined APIs and GCP tools:
 - Identify various **fields within a given claims form**
 - Identify handwriting and **transcribe into a flat format (.csv)**

GCP Capabilities Exploration

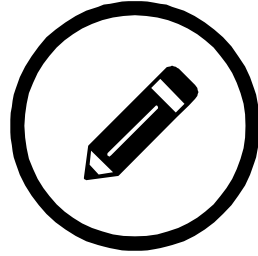


Plan

Review API Documentation

Identify API pricing, documentation, and samples.

The underlying models from these selections will be used to drive the remainder to the exploration.

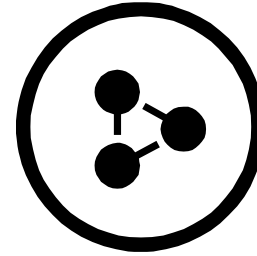


Design

Gather Training Data

Choose several forms of interest and create a training sample to test the output of the selected models.

Having a variety of forms will help understand the limitations of pre-built solutions by determining common areas of weakness.

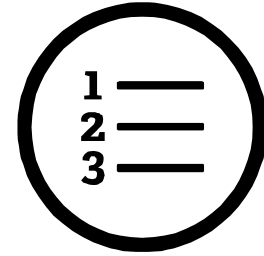


Build

Develop and Test Models

Develop a sample pipeline and test APIs across a variety of criteria including accuracy and compute time.

Outputting the results to a flat format allows for analysis and standardization across forms.



Evaluate

Review Results & Accuracy

Using the criteria outlined in the program objectives, determine whether to continue research into GCP APIs.

Socialize results and suggest areas for future work.



Cloud Vision API Pricing

Applicable	Feature Type	Description	Price per 1000 Units	
			Units 1001 - 5,000,000 / month	Units 5,000,001 - 20,000,000 / month
⊗	LABEL_DETECTION	Add labels based on image content (see Detecting Labels)	\$1.50	1.00
✓	TEXT_DETECTION	Perform Optical Character Recognition (OCR) on text within the image (see Detecting Text)	\$1.50	\$0.60
✓	DOCUMENT_TEXT_DETECTION	Perform OCR on dense text images, such as documents (see Document Text Detection)	\$1.50	\$0.60
⊗	SAFE_SEARCH_DETECTION	Determine image safe search properties on the image (see Detecting Safe Search Properties)	Free with Label Detection, or \$1.50	Free with Label Detection, or \$0.60
⊗	FACE_DETECTION	Detect faces within the image (see Detecting Faces)	\$1.50	\$0.60
⊗	LANDMARK_DETECTION	Detect geographic landmarks within the image (see Detecting Landmarks)	\$1.50	\$0.60
✓	LOGO_DETECTION	Detect company logos within the image (see Detecting Logos)	\$1.50	\$0.60
⊗	IMAGE_PROPERTIES	Compute a set of properties about the image, such as the image's dominant colors (see Detecting Image Properties)	\$1.50	\$0.60
⊗	WEB_DETECTION	Detect topical entities such as news, events, or celebrities within the image	\$3.50	Contact Google for more information
✓	CROP_HINTS	Determine suggested vertices for a crop region on an image (see Detecting Crop Hints)	Free with Image Properties, or \$1.50	Free with Image Properties, or \$0.60
⊗	OBJECT_LOCALIZATION	Detect and extract multiple objects in an image (see Detect Multiple Objects)	\$2.25	\$1.50



Required Capabilities

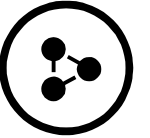
A successful extraction system needs to accurately transcribe handwritten content and associate it with the correct field

The image shows a scanned 'Member Claim Form' from Anthem. The form is annotated with numbered circles 1 through 6, indicating specific capabilities required for extraction:

- 1: Points to the 'Anthem' logo at the top right.
- 2: Points to the 'Member Claim Form' title at the top left.
- 3: Points to the 'Section A. PATIENT INFORMATION' header.
- 4: Points to a handwritten '10/10' in the 'Date of birth' field.
- 5: Points to the 'Section B. SUBSCRIBER INFORMATION (on Anthem Blue Cross card)' header.
- 6: Points to the 'Section C. MEDICAL INFORMATION' header.

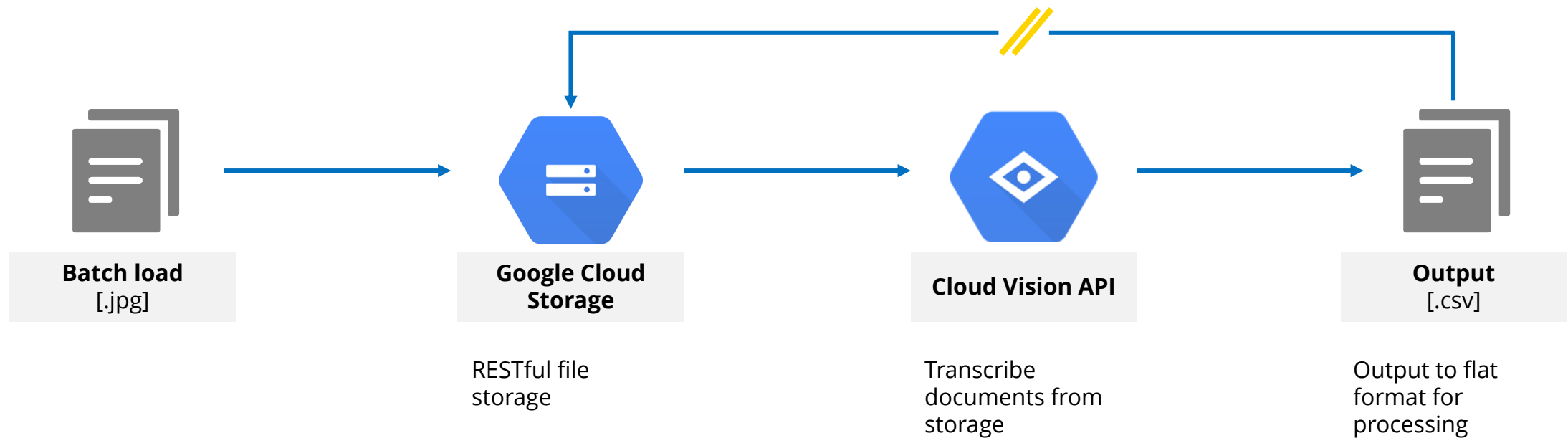
#	Capability	Vision API
1	Correct skewness, remove noise, and crop	✓ ¹
2	Detect form template	✓
3	Identify relevant section content	✓
4	Transcribe handwritten text	✓
5	Standardize, validate, and extract content from data fields	X ²
6	Identify checked boxes and other non-text symbols	✓

¹ Crop Hints (Beta)
² Data Validation feature not available in Vision API



Proof of Concept Architecture

A basic workflow to test multiple claim types with the ability to scale using products like Google Dataflow and BigQuery



Represents future additions

- [Dataflow](#) for batch or streaming data transformations and processing
- [BigQuery](#) for inexpensive database storage with interactive analysis tools

Claim Forms

Multiple claim forms were compared to evaluate the performance of the Vision API across different document types

New York Prescription Form

Prescription Drug Reimbursement / Coordination of Benefits Claim Form
An incomplete form may delay your reimbursement.
See the back for instructions and complete all information.

Cardholder Information See your prescription drug ID card.
Group No.
Member ID:
Member Name First Last
Street Address
City State ZIP
City State ZIP

Patient Information
Patient Name First Last
Patient Date of Birth (Month/Day/Year)
Sex Relationship to Patient Member
☐ Female ☐ Male ☐ 1 Self ☐ 2 Spouse ☐ 3 Eligible Child ☐ 4 Dependent Student ☐ 5 Disabled Dependent ☐ 6 Dependent Parent ☐ 7 Non-spouse Partner ☐ 8 Other

Pharmacy Information
Name of Pharmacy
Street Address
City State ZIP
Telephone (include area code)
Is this an on-site nursing home pharmacy? ☐ Yes ☐ No
Signature of Member Date
Signature of Pharmacy Date

Claim Receipts
Type receipts in shaded cells in the back.
Check the appropriate box if any receipts or bills are for a:
☐ Compound prescription
Make sure your pharmacist fills ALL the VALID RUC numbers, cost and quantities for each ingredient on the back of this form and attach receipts. Claims will be returned if incomplete.
ONE CLAIM FORM PER COMPOUND SUBMISSION
☐ Medication purchased outside of the United States
Please include date.
Country
Currency used
☐ Allergy medication
Coordination of Benefits
If another health plan has paid a portion, mark the appropriate box for your primary coverage method. See the back for more information.
Is this a coordination of benefits claim?
☐ Yes ☐ No
Another health plan and one you are enclosing a statement that both has been paid and how much the other plan or paid (1)
☐ Card Payment (2)
☐ Domestic Scripts Mail Order (3)
Any person who is knowingly and with intent to defraud, or cause of damage any insurance company, submit a claim or receipt containing any information (false, deceptive, incomplete, or misleading) when they apply for a health claim may be committing a fraudulent insurance act, which is a crime and the subject of a criminal action. In addition, penalties, including fines and/or imprisonment or death of benefits.
Please tape receipts on the back of this page.

California Medical Claim Form

Member Claim Form
Anthem
Please use a separate claim form for each patient. Your cooperation in completing all items on the claim form and attaching all required documentation will help expedite quick and accurate processing. SEE REVERSE SIDE FOR COMPLETE INSTRUCTIONS.

Section A: PATIENT INFORMATION
Last Name First Name MI
Does the patient have other health insurance coverage? ☐ No ☐ Yes
☐ Self ☐ Spouse ☐ Son ☐ Daughter ☐ Other
Name of other health insurance company Group No. Employee Name Policy No.
Date of Birth (MM/DD/YYYY)

Section B: SUBSCRIBER INFORMATION (On Anthem Blue Cross card)
Last Name First Name MI
Street Address (please include apt. no.)
City State ZIP Code
Home phone no. Work phone no. Date of Birth (MM/DD/YYYY)

Section C: MEDICAL INFORMATION
HEALTH CARE PROVIDER: Due to this section to report any HEALTH CARE provider that has not already been reported in this form. Use this section for the physician, clinic, ambulance company, private duty nurse, etc. Attach itemized bill or pharmacy. Please be sure to duplicate bills for each service.
Was this portion expense the result of an accident? ☐ Yes ☐ No
Was this condition or injury pre-existing? ☐ Yes ☐ No
Have you filed for Workers' Compensation? ☐ Yes ☐ No
When did this injury or accident occur? (MM/DD/YYYY)
Diagnosis code Procedure code Unit

RULES MUST BE FOLLOWED
Completed forms with receipts and/or itemized "balance bill" statements cannot be processed. Each itemized bill must include:
• Name and address of provider (doctor, hospital, laboratory, ambulance service, etc.)
• Name of patient
• Service provided
• Date of service
• Amount charged for each service
• Diagnosis code
• Procedure code
• Date of bill
Identify that, to the best of your knowledge, the information on this Member Claim Form is true and correct. I authorize the release of any and all information necessary to process this claim.
Signature Name Date

New York Medical Claim Form

HEALTH INSURANCE CLAIM FORM MEMBER SUBMITTED
FOR 3407 CHURCH STREET STATION, NEW YORK, NY 10008-3407
APPROVED OMB-0138-0038
NOTE: Important filing instructions on next page.

PATIENT AND ASSURED INFORMATION
Last Name First Name MI
Date of Birth (MM/DD/YYYY)
Sex ☐ Male ☐ Female
Address
City State ZIP
Home phone Work phone
Signature of Member Date
Signature of Provider Date

PATIENT AND ASSURED INFORMATION
Last Name First Name MI
Date of Birth (MM/DD/YYYY)
Sex ☐ Male ☐ Female
Address
City State ZIP
Home phone Work phone
Signature of Member Date
Signature of Provider Date

PATIENT AND ASSURED INFORMATION
Last Name First Name MI
Date of Birth (MM/DD/YYYY)
Sex ☐ Male ☐ Female
Address
City State ZIP
Home phone Work phone
Signature of Member Date
Signature of Provider Date

- 25 total fields (20 text / 5 symbolic)
- Text and symbolic entry
- **Poses additional difficulty for Vision API due to bounded characters**

- 29 total fields (23 text / 6 symbolic)
- Bounded text boxes
- Free form text entry
- **Similar character bounding difficulty**

- 60 total fields (47 text / 13 symbolic)
- Multiple rows / columns
- Numeric and text free form
- **Overlap of handwriting in freeform area**

Entry Types

Several entry types were used to understand how model performance differs within identical forms

Blank

Form for each patient. Your use of this form is subject to the terms and conditions of the license agreement. For more information, see the back and accurate processing. SEE REVERSE.

PATIENT INFORMATION

Does the patient have other health insurance coverage? ☐ Yes ☐ No

Relation to subscriber: ☐ Self ☐ Spouse ☐ Other

Name of other health insurance company: _____ Group no.: _____ Employer: _____

Section B. SUBSCRIBER INFORMATION (on Anthem Blue Cross card)

Identification no.: _____ Group: _____

Last name: _____ First name: _____

Address (please include apt. no.): _____

- Identifies bounds for fields
- Baseline for model performance
- Useful for form identification

Typed

Form for each patient. Your use of this form is subject to the terms and conditions of the license agreement. For more information, see the back and accurate processing. SEE REVERSE.

PATIENT INFORMATION

Does the patient have other health insurance coverage? ☒ Yes ☐ No

Relation to subscriber: ☒ Self ☐ Spouse ☐ Other

Name of other health insurance company: Anthem Group no.: 1234 Employer: ABC

Section B. SUBSCRIBER INFORMATION (on Anthem Blue Cross card)

Identification no.: 01234 Group: 567

Last name: T R E I D E S First name: L

Address (please include apt. no.): 3 MAIN ST

- Variable performance across forms
- Affected by distance and spacing between fields

Handwritten

Form for each patient. Your use of this form is subject to the terms and conditions of the license agreement. For more information, see the back and accurate processing. SEE REVERSE.

PATIENT INFORMATION

Does the patient have other health insurance coverage? ☒ Yes ☐ No

Relation to subscriber: ☒ Self ☐ Spouse ☐ Other

Name of other health insurance company: ANTHEM Group no.: 1234 Employer: AB

Section B. SUBSCRIBER INFORMATION (on Anthem Blue Cross card)

Identification no.: 01234 Group: _____

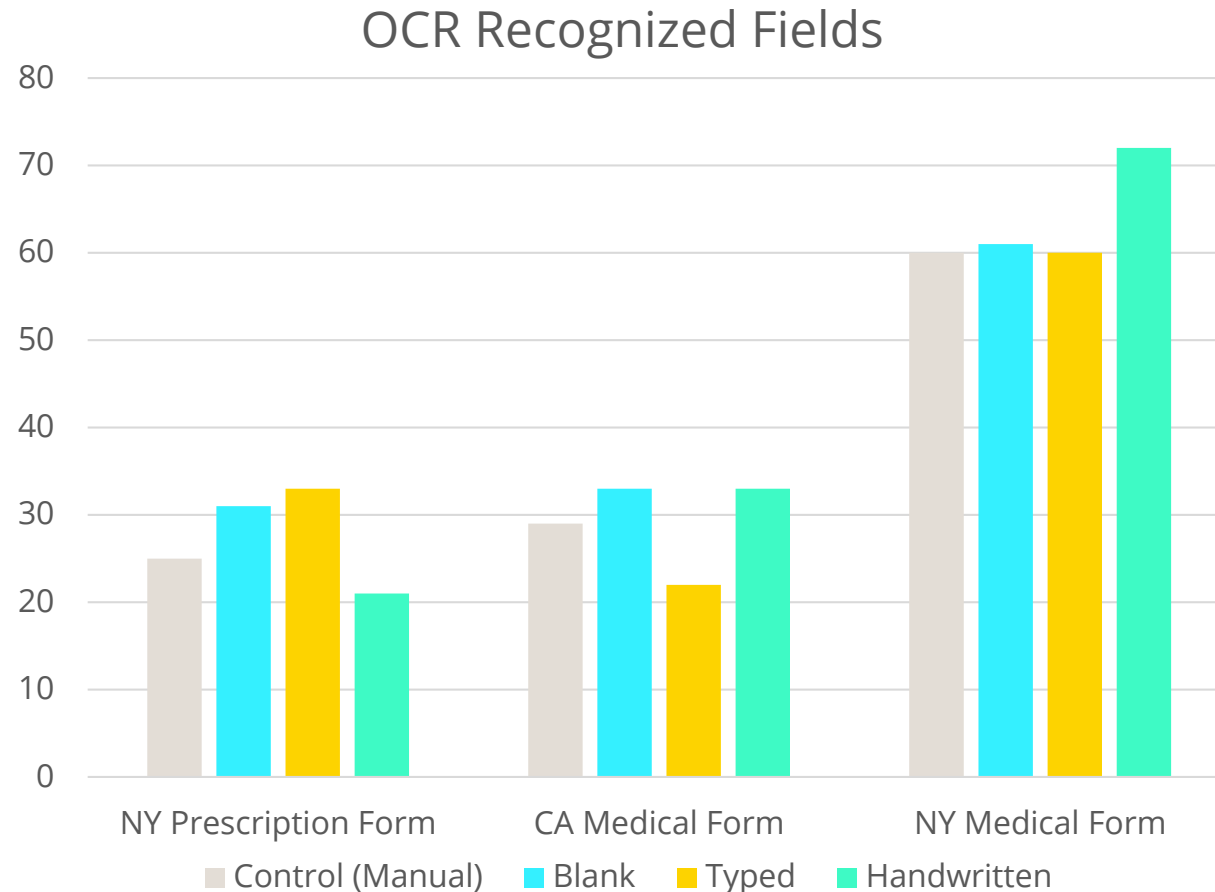
Last name: T R E I D E S First name: L

Address (please include apt. no.): 3 MAIN ST

- Highly variable performance within forms
- Affected by a variety of factors (i.e. pen color, size, etc)

OCR: Segmentation

Testing the API's auto segmentation feature



API Segmentation is Highly Variable

Even in forms where the amount of fields is relatively similar, the content included in each segment is inconsistent and/or overlapping

Not All Information is Captured

When automatic segmentation is performed by the API, some information is omitted or truncated due to being outside of the bounds

Preprocessing suggested to define field boundaries for each form type

OCR: Accuracy

Isolated fields were compared to determine accuracy of handwritten vs. typed text

Field Name	Data Type	Actual	Handwritten	Typed
Last Name	Text	ATREIDES	ATRIEL PIESE	ALTRE DIES
Group Number	Numeric	56789	56789	56789
Coverage*	Symbolic (Yes / No)	No	N/A	N/A
Subscriber Relation*	Symbolic (Multiple)	Self	N/A	N/A

Last name

A | T | R | E | I | D | E | S | |

Group no.

56789

Does the patient have other health insurance coverage?

☐ Yes ☒ No

Relation to subscriber

☒ Self ☐ Spouse ☐ Son ☐ Daughter

*Symbols to be handled by AutoML preprocessing



Suggestions for Future Work

Addition of pre-processing and form validation would dramatically improve ability to discern individual fields



Segmentation Based on Form Type

- Removes variance within forms caused by font placement, size, field overlap
- Requires additional API calls
 - API Call #1: Form recognition
 - API Call #2: OCR

Data Validation

- Preprocessing includes data validation rules to guide field specific values and recognize checkboxes
- Suggested validation rules:
 - Data value (alphanumeric, symbol)
 - Data format (dd/mm/yy, phone #)

Thank you.

Philip Mohun

Deloitte Consulting LLP

GCP Professional Data Engineer (*Candidate*)

Contact: phmohun@deloitte.com

This publication contains general information only, and none of the member firms of Deloitte Touche Tohmatsu Limited, its member firms, or their related entities (collective, the “Deloitte Network”) is, by means of this publication, rendering professional advice or services. Before making any decision or taking any action that may affect your business, you should consult a qualified professional adviser. No entity in the Deloitte Network shall be responsible for any loss whatsoever sustained by any person who relies on this publication.

As used in this document, “Deloitte” means Deloitte Consulting LLP, a subsidiary of Deloitte LLP. Please see www.deloitte.com/us/about for a detailed description of the legal structure of Deloitte USA LLP, Deloitte LLP and their respective subsidiaries. Certain services may not be available to attest clients under the rules and regulations of public accounting.

**Copyright © 2018 Deloitte Development LLC.
All rights reserved. Member of Deloitte Touche Tohmatsu Limited**