



# Truly refreshing document analysis



MINT.extract

Automated Data Extraction from Documents  
using Machine Learning

19.06.2019

Swisstext 2019

Aaron Richiger, Co-Founder turicode

# Agenda

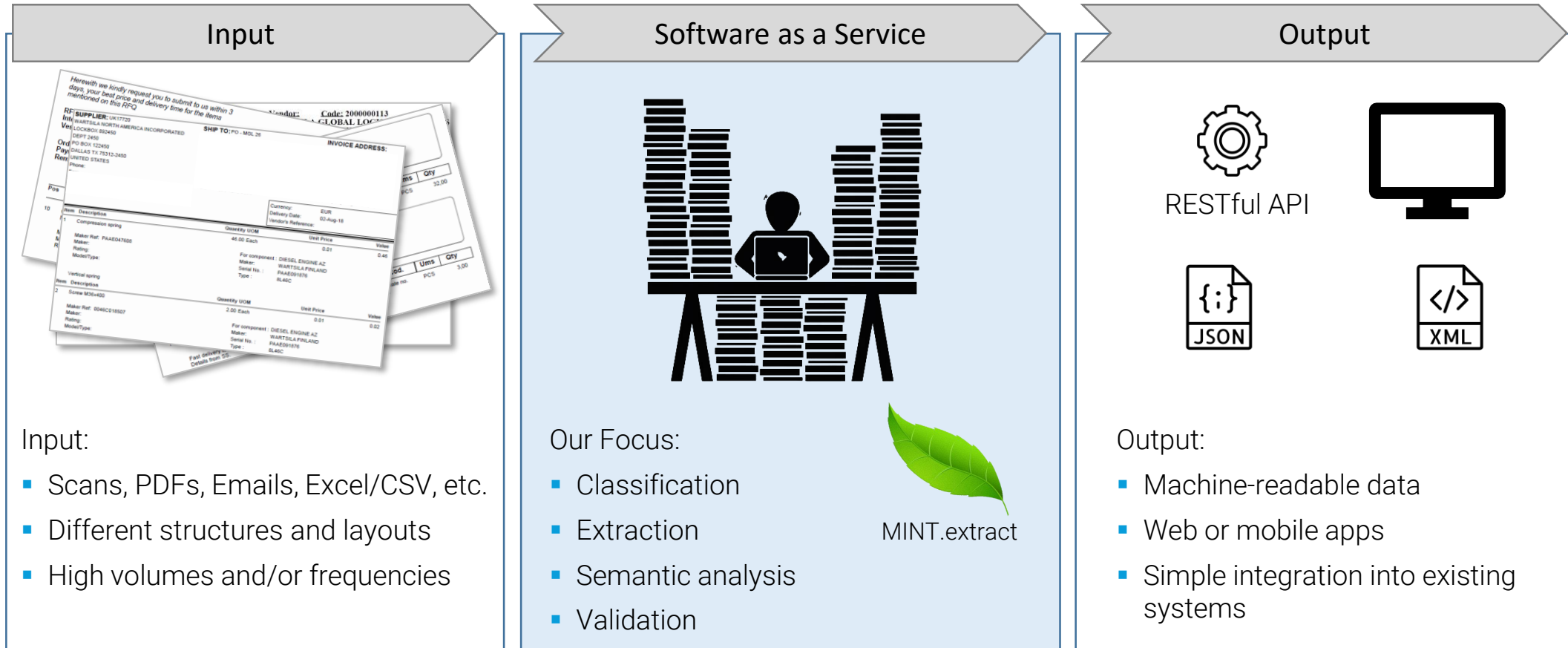


## MINT.extract

- What it is
- How it works
- Results
- How to use it



# Problem Statement



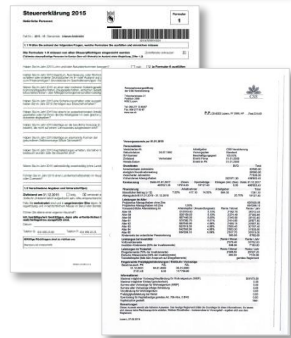
# How MINT.extract structures your documents



MINT.extract



OCR / Text Recognition



Pension Statements



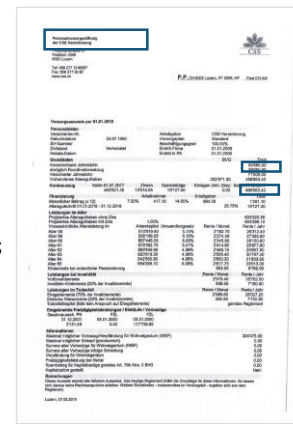
Tax Statements



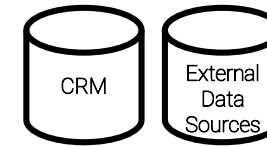
Bank Statements



Insurance Policies



Extraction of text, images etc. with Machine Learning



Automated Validation and Enrichment



RESTful API



```
{  
  "Pensionfund":  
  "Personalvorsorgestiftung der CSS Versicherung",  
  "Salary": 99585 ,  
  "Saldo": 496863  
}
```

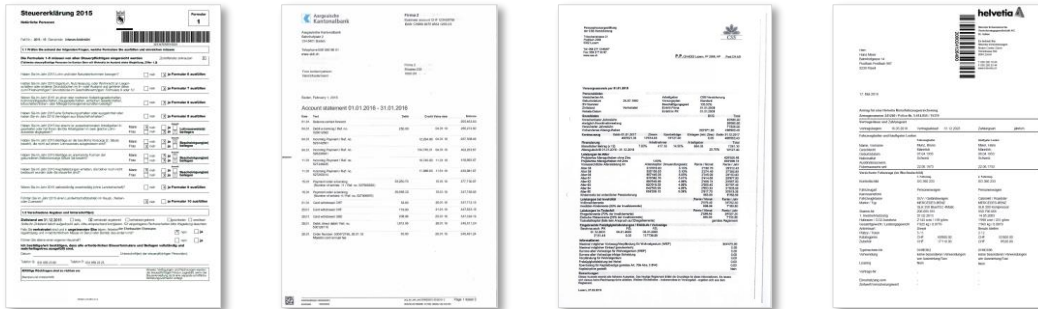


Web-Application for visual inspection



# How it works: Classification

## Image Based Classification



## Text Based Classification

Gerüchte machten schon eine Weile die Runde, nun werden sie wahr: Facebook kündigt die Einführung einer Kryptowährung namens Libra an. Nicht sofort, sondern Anfang bis Mitte nächsten Jahres. Das Ziel ist ambitioniert. So soll sich die Libra als erste »Weltwährung« überaupt etablieren und zusammen mit der Technologie, auf der sie basiert, das Leben von Milliarden von Menschen erleichtern. Das klingt ausserst anspruchsvoll und fast wie ein Märchen – aber was ist dran?

Die Überlegungen gehen zunächst von der Beobachtung aus, dass man heute aufgrund des Internets und moderner Kommunikationmöglichkeiten praktisch nur noch ein günstiges Smartphone benötigt, um sich zu informieren, um Nachrichten beinahe kostenlos auszutauschen oder um sich den Alltag mit der Inanspruchnahme allerlei günstiger Dienstleistungen hoher Qualität zu versüssen. Die Technologieunternehmen haben ganze Arbeit geleistet, indem sie den Konsumenten den Zugang zu traditionellen Warenangeboten und Dienstleistungen vereinfacht und vergünstigt haben oder indem sie neue einführen.

Der amerikanische Internetkonzern macht Genf zum Zentrum seiner Zukunftsvision vom »Internet of Money«. Die neue Kryptowährung Libra soll den Umgang mit Geld so einfach machen wie das Verschieben einer SMS. Hat das Projekt Erfolg, könnte es das Finanzsystem umkrempeln. Der amerikanische Internetkonzern macht Genf zum Zentrum seiner Zukunftsvision vom »Internet of Money«.

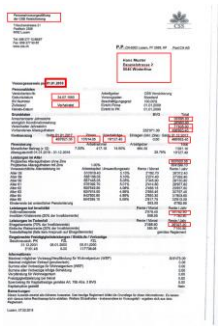


## Combination of Text and Image Based Classification

- Showed the best results
- Importance of text or image input depends on the documents



# How it works: Data Extraction



1. Word level classification based on
  - Text value
  - Textual context
  - Layout information



2. Grouping datapoints to structured entities

17\$

3. Character level classification



# How it works: Data Extraction

Quick Extract Training Production Train F1: 0.91

Production > Email\_Auftragsbestaetigung\_100577670\_ORD\_1054040925.pdf

Page 1 of 2 Zoom: [ ]

Address: **Herr Martin Keller, Plattenstrasse 34, CH-8032 Zurich**

CustomerNo: **1832250**  
 Reference: **1054040925**  
 OrderDate: **15.12.2016**

Ihre Kundennummer: **1832250**  
 Ihre Auftragsnummer: **1054040925**  
 Auftragsdatum: **15.12.2016**  
 Zahlung des Auftrages: **Kreditkarte**

Sehr geehrter Herr Keller

Vielen Dank für Ihren Auftrag bei buch.ch.  
 Gerne bestätigen wir den Eingang Ihres Auftrags:

| Medium | Menge | Autor und Titel                                           | ISBN / EAN    | Betrag           |
|--------|-------|-----------------------------------------------------------|---------------|------------------|
| Buch   | 1     | Bernd Regenberg / Regenberg, B: Jahrbuch Lastwagen 2017   | 9783861338178 | 21.90 CHF        |
| Buch   | 1     | Felicity Brooks / Brooks, F: Meine kleine Stickerwelt: La | 9781782323167 | 9.40 CHF         |
| Buch   | 1     | Haruki Murakami / Murakami, H: Kafka am Strand            | 9783442733231 | 17.90 CHF        |
| Buch   | 1     | Haruki Murakami / Murakami, H: Von Männern, die keine F   | 9783442714254 | 14.90 CHF        |
| Buch   | 1     | Peter Frankopan / Frankopan, P: Silk Roads                | 9781410883997 | 19.90 CHF        |
| Buch   | 1     | W. Timothy Gallwey / Gallwey, W: The Inner Game of Tennis | 9781447288503 | 19.90 CHF        |
|        |       | Gutscheinbetrag                                           |               | - 30.00 CHF      |
|        |       | Versandkosten                                             |               | 0.00 CHF         |
|        |       | <b>Total</b>                                              |               | <b>73.90 CHF</b> |

Address: Herr Martin Keller Plattenstrasse 34 CH 8032 Zurich

CustomerNo: 1832250

Reference: 1054040925

OrderDate: 15.12.2016

Items

- 
- Amount: 1

Description: Felicity Brooks / Brooks, F: Meine kleine Stickerwelt: La

ProductNo: 9781782323167

Status:

Price: 9.4

▲ Must be >= 0.5

Currency:
- 
- 
- 
-

# How it works: Analytics & Validation

Based on the extracted data, we build project-specific ML systems for further analytics and validation:

## Data Analytics with ML

- NLP (e.g. NER, topic detection, etc.)
- Predictive Maintenance
- Trigger a specific workflow



## Validation with ML

- Anomaly detection
- Validation classification (critical vs. non-critical validation violations)







# Results: Classification

Documents: Financial documents  
Case: Incoming mail classification  
Data: 17k sample documents for 4 classes  
Approach: Combined text- and image-based classification (80/20 split with CV)

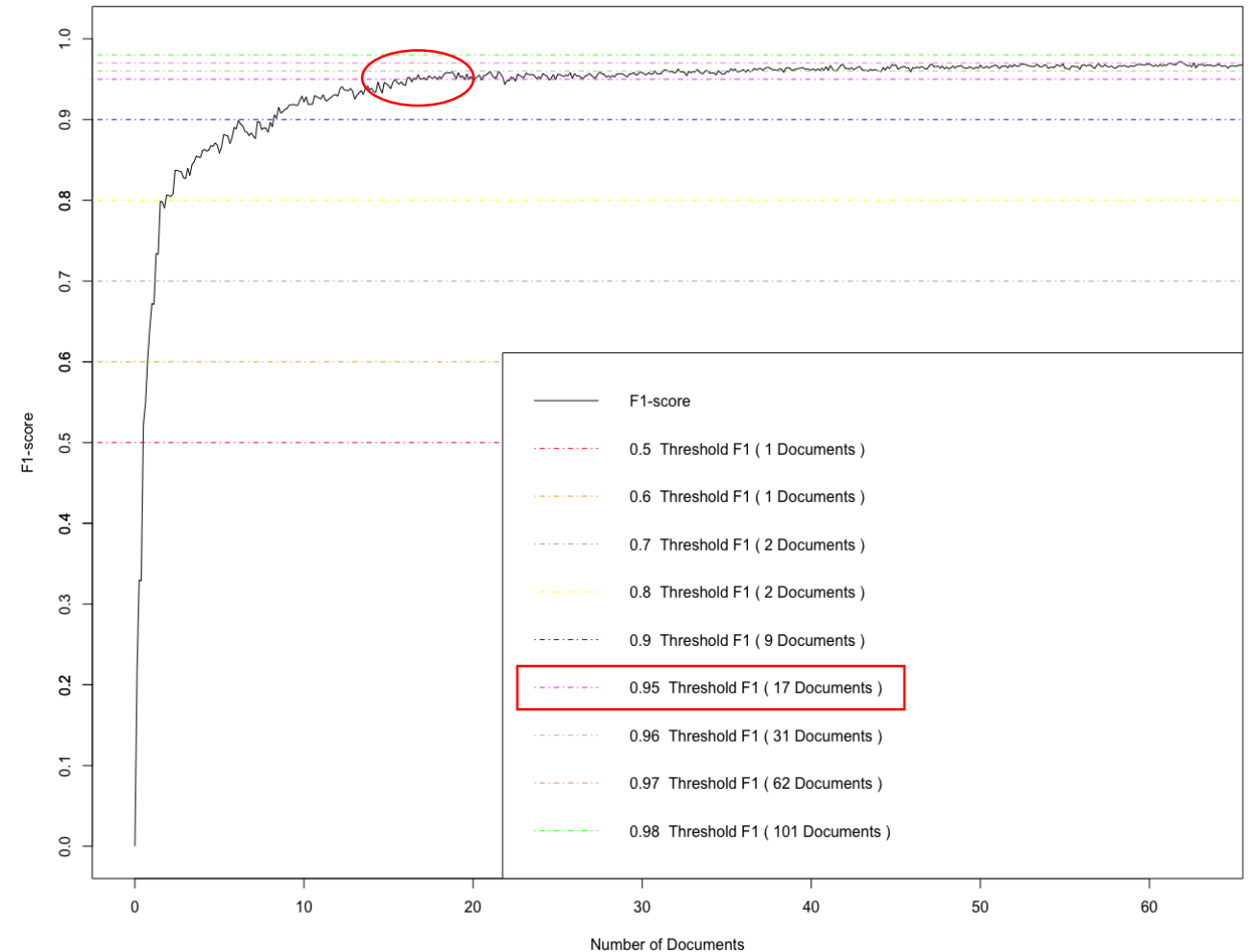
| Class                                   | Precision   | Recall      | F1 score    |
|-----------------------------------------|-------------|-------------|-------------|
| Official Documents (Amtliche Dokumente) | 0.98        | 0.99        | 0.99        |
| Card Application (Kartenanträge)        | 0.99        | 0.99        | 0.99        |
| Direct Debiting Application (LSV)       | 0.99        | 0.99        | 0.99        |
| Payment Order (Zahlungsauftrag)         | 0.96        | 0.98        | 0.97        |
| <b>Weighted Average</b>                 | <b>0.99</b> | <b>0.99</b> | <b>0.99</b> |

# Results: Data Extraction



MINT.extract

Documents: Purchase Orders  
Case: Data extraction and validation  
Data: 70 documents from 3 customers with 8 data points to extract  
Approach: Word level classification with 80/20 split and CV





# How to use MINT.extract

## Option 1: Service provided by turicode

- You send us the documents and the spec
- We provide you with:
  - REST API (for ongoing services)
  - Excel, JSON, XML, TXT (for projects)



## Option 2: Try and train it yourself

- Request access to MINT.extract on our website [turicode.com](https://turicode.com)
- We provide you access to your own MINT.extract web application
- Upload, annotate and process your documents



# Excerpt from our References



MINT.extract



Gesellschaft für Schweizerische  
Kunstgeschichte



Stadt Zürich



Universität  
Zürich<sup>UZH</sup>

Supported by:



Co-funded by the Horizon 2020 programme  
of the European Union

# Contact



turicode AG  
Technoparkstrasse 2  
CH-8406 Winterthur

Aaron Richiger

aaron.richiger@turicode.ch

+41 79 613 71 78

[www.turicode.com](http://www.turicode.com)



swiss made software