

Ockr Specification

Bormann, Fabian
fabian.bormann@ockr.io



November 26, 2023

Contents

1	Abstract	2
2	Motivation	3
3	Process	4
3.1	Document creation	4
3.1.1	Information extraction	4
3.2	OCR verification	4
3.3	Guarantee authenticity	4
3.4	Authorship verification	4
3.5	Trust Trees	5
4	Algorithms	6
4.1	OCR Algorithms	6
4.2	Hash Algorithms	6
4.3	Puzzle Algorithms	6
4.3.1	Default Puzzle Algorithm	6
4.4	Templates	7
5	Implementation	9
5.1	Ockr REST API	9
5.1.1	Endpoints	9
5.2	Ockr Frontend	10
5.3	Ockr Model Zoo	11
5.4	Ockr OCR containers	11
5.5	Ockr plugins	11
5.6	Ockr App	11
5.7	Chain Submitter Service	11

1 Abstract

Ockr aims to set a standard for creating machine-readable and reliable documents, enabling the verification of their authenticity. In short Ockr tries to tackle three problems:

1. How can we make sure that a document is machine-readable?
2. How can we ensure that a document has not been modified?
3. How can we ensure that a document really has been issued by a certain authority?

Process automation is an important topic across different industries as well as certification processes or even for governments. In the past years many processes became paperless and digital. However, there are still processes that following hybrid approaches or are still fully paper-based. If many different parties are involved in a process, it is often hard to keep track of the current status for example in logistics, medicine, student certifications or documents issued by the government. Blockchain technology can be used as a trust anchor to ensure the authenticity of a document. However, the information itself is not stored on-chain. Files are used to store the information and hashes of these files are stored on-chain. Due to compression techniques or printing the file remains visually the same but the hash changes. In logistics it is often required to attach a paper document to the actual good because it is not possible to share files with all parties involved. Many Universities and schools are issuing digital certificates nowadays but to apply for a job, the certificate needs sometimes compression techniques or even printing to be shared with the employer. There are many different scenarios having the same problem: The document has been created digitally but due to printing, compression or conversions, it is not machine-readable anymore. Ockr aims to solve this problem by providing a standard for creating machine-readable and reliable documents. A QR code is used to store a hash of the actual information alongside the document. The QR code also contains metadata supporting an OCR algorithm to re-create the correct information. This ensures reliable process automation in hybrid environments and decouples the actual information from the file container.

2 Motivation

Most of the documents today are created in a digital form and are shared online. The file itself serves as a vehicle for the actual information. File hashes are often used to prove the authenticity of the information. However, the file hash can be changed over time, while the information remains the same. Compression techniques are used behind the scenes while sharing the files online, as well as storing the file in cloud environments. Finally printing the file and scanning it again will result in a different file hash. Images or documents that store the information in images are not machine readable. While process automation became crucial in many industries, many processes still rely on manual work because of the lack of machine readable documents. OCR (Optical Character Recognition) is a technique that can be used to extract the information from images, but it is not reliable enough to be used in critical processes as it can not guarantee the authenticity of the information. Ockr aims to solve these problems by providing a standard for creating machine-readable and reliable information.

3 Process

To decouple the information from the file, the raw information needs to be provided and a hash of the information needs to be calculated. This hash will be used to verify the authenticity of the information and stored within a QR code on the document. Besides the hash, the QR code will also contain metadata helping a machine to re-create the information.

3.1 Document creation

As most of the documents are created digitally, the raw text can be extracted from the file such as DOCX, XLSX, PDF and other formats. In case there is no raw text available, OCR can be used to extract the text from the document but a human would need to verify the text.

3.1.1 Information extraction

To ensure a good user experience, plugins for common document editors will be provided to extract the raw information from the document and to create the Ockr QR code. The plugins will be open source and can be used as a reference implementation for other plugins.

3.2 OCR verification

There is no internet connection needed and no additional service required to verify the OCR result. Therefore data pipelines can ensure to have a reliable OCR result in critical processes. An algorithm breaks down the document in different parts and verifies the OCR result of each part. The algorithm can be implemented in any programming language and can be used in any environment. Besides the main hash of the information a hash of each part will be calculated and stored within the QR code. There are chars that are often misinterpreted by OCR like 0 and O. With having the hash of each part, different combinations of chars can be tried to find the correct result. Algorithms and OCR models will be provided as open source projects under the ockr.io umbrella.

3.3 Guarantee authenticity

An open source backend service will be publicly available to register an information hash on a blockchain. The service will be available as a REST API. Once the information hash has been settled on-chain, it can be verified by anyone. The service will be open source and can be used as a reference implementation for other services.

3.4 Authorship verification

The authority that registered the information would also need a registration event on-chain. Anyone can register a new authority but to give a certain

authority more credibility, the authority can be verified by other authorities.

3.5 Trust Trees

The trust tree is a graph of authorities that trust each other. An algorithm will be provided to calculate the trust tree and its score. The score can be used to determine the credibility of an authority. The REST API will also offer an endpoint to check if a certain authority is part of the trust tree of another authority. If for example, an automated HR process wants to verify a certificate, it can check if the university that issued the certificate is trusted by the government.

4 Algorithms

Algorithms need to ensure that an OCR result is reliable. The first set of algorithms will be used to extract the information from the document. Besides of the OCR algorithm itself, they will also include algorithms to build the information hash from the OCR output. Those will be open source and are part of the Ockr umbrella. They provide something like the syntax for the OCR scan results. The second set of algorithms will be used to give the results a semantic. Those are called templates and they are optional. Templates will not be part of any repository. They can be registered on-chain using the REST API. Templates can be used to give an automated process the relevant values. For example a value relative to certain anchor words could be defined as the product weight on a delivery note to calculate the shipping costs for an automated process.

4.1 OCR Algorithms

PP-OCRv4-mobile and PP-OCRv4-server are part of PaddleOCR and will be provided as ONNX models in the Ockr model zoo. The model zoo will be an open source repository where everyone can contribute OCR models or enhance the existing ones. A model should run isolated and wrapped by a REST API encapsulated in a Docker container. The blueprint will use Fast API and the Docker image will be provided alongside to the code as a reference implementation.

4.2 Hash Algorithms

SHA256 will be used as the hash algorithm. The hash will be calculated from the initial information during the content creation process. In the verification process it will be calculated from the OCR result and a naive string concatenation approach. As this solution is not reliable enough, SHA256 will be applied also to the output of a post-processing method, introduced as puzzle algorithms.

4.3 Puzzle Algorithms

The puzzle algorithms will be used to break down the image in different parts. Each part will be hashed by SHA256 and stored within the QR code. The puzzle algorithms will be open source and can be used as a reference implementation for other algorithms. A simple and default algorithm will be described in the following section.

4.3.1 Default Puzzle Algorithm

The default algorithm takes the image and splits it into n parts. It takes x and y as parameters to define the number of parts. The image will be split into x parts horizontally and y parts vertically. The parts will be hashed by SHA256 and stored within the QR code.

Algorithm 1 Default Puzzle Algorithm

Require: *OcrResults***Require:** $x \geq 1$ **Require:** $y \geq 1$

```
subHashes  $\leftarrow$  []  
for  $i \leftarrow x$  to  $x$  do  
  for  $j \leftarrow y$  to  $y$  do  
    results  $\leftarrow$  getOcrResultsUnderArea( $x, y$ )  
    hash  $\leftarrow$  SHA256(concat(results))  
    subHashes  $\leftarrow$  append(hash)  
  end for  
end for  
return subHashes
```

4.4 Templates

A template consists of a set of anchors, such as the OR code or other required words, and a set of keys including relative positions to a subset of anchors. The anchors are used to find the values for the keys in the OCR result. The values can be defined as a match to a key, or as a regex matching the key. The template structure will be defined as a JSON schema and can be registered on-chain using the REST API and a Cardano metadata standard. An offline template editor will be provided to create and edit templates. The editor will be open source and can be used as a reference implementation for other validators or editors.

```
interface Template {  
  anchors: Anchor[];  
  keys: Key[];  
}  
  
interface Anchor {  
  type: AnchorType;  
  reference: BorderSide | string;  
}  
  
enum BorderSide {  
  Top,  
  Right,  
  Bottom,  
  Left  
}  
  
enum AnchorType {  
  QRCode,  
  Text,  
  Border
```

```
}

// regex is optional otherwise the full ocr result will be used
interface Key {
    name: string;
    regex: string;
    matcher: {
        anchor: number;
        position: Position;
    }
}

// relative to page width between -1 and 1
interface Position {
    x: number;
    y: number;
    tolerance: number;
}
}
```

5 Implementation

There are different parts mentioned in the algorithm and process section. This section will describe the implementation of those parts and the concret implementation.

5.1 Ockr REST API

The Ockr REST API is the main interface to interact with the Ockr ecosystem. It is used to register information hashes on-chain, to register authorities and to register templates. The API is implemented using Java Sprint Boot. It will be available as an open source project under the Ockr umbrella. The API will be available as a Docker image which allows companies to run it in their own environment instead of using the public API. The public API will be available at *api.ockr.io/v1*.

5.1.1 Endpoints

```
// Registers a new information hash on-chain

POST /register/hash

// Request Body
{
  "hash": "string",
  "signature": "string",
  "template": "string" | undefined,
  "subHashes": ["string"] | undefined,
  "algorithm": "string" | undefined,
}

// Registers a new template on-chain
// Template is a JSON schema described in the algorithms section
// The template will be minified and validated

POST /register/template

// Request Body
{
  "template": Template,
  "signature": "string",
}

// Registers a new authority on-chain
// The authority can be co-signed by the registering service as the
  first validator

POST /register/authority
```

```

// Request Body
{
  "name": "string",
  "signature": "string",
}

// Returns the trust tree of an authority if the hash matches with a
  registered information hash on-chain

POST /verify/hash

// Request Body
{
  "hash": "string",
  "template": "string" | undefined,
  "subHashes": ["string"] | undefined,
}

// Returns all templates that have been registered on-chain

GET /template

// Returns a template by id

GET /template/{id}

// Returns all models connected to the api

GET /model

// Connects a model to the api

POST /connect/model

```

The API will implement a yaci store starter to fetch templates, authorities and information hashes from the blockchain. Those will be stored in a local database.

5.2 Ockr Frontend

The Ockr frontend is a web application that can be used to validate and register information hashes on-chain, to register and look up authorities and to register and test templates. The frontend is implemented using React and will be available as an open source project under the Ockr umbrella. The frontend will be available at *ockr.io*.

5.3 Ockr Model Zoo

The Ockr model zoo is a collection of OCR models that can be used to extract information from documents. The preferred model format is ONNX. The model zoo will be available as an open source project under the Ockr umbrella.

5.4 Ockr OCR containers

An Ockr OCR container offers a FastAPI REST-service within in a docker container, that can be used to extract information from documents. It uses models e.g. from the Ockr model zoo to extract the information. The containers will be available as an open source projects under the Ockr umbrella. In its environment variables the container has a reference to the Ockr REST API. The container will connect itself to the Ockr REST API on startup. The API can be configured to block new connections. Besides a health check endpoint, the container offers an endpoint to perform the actual OCR. The OCR endpoint will be available at *localhost:8000/predict [POST]*. It assumes a base64 encoded image. The output will be a JSON object containing the OCR result.

5.5 Ockr plugins

Ockr plugins are plugins for common document editors that can be used to extract the information from the document and to create the Ockr QR code. The plugins will be available as an open source projects under the Ockr umbrella. The first plugins/add-ons on the roadmap are for Google Docs, Google Sheet, Google Slides, Microsoft Word, Microsoft Powerpoint, Microsoft Excel. The first iteration of the plugins will only provide a QR code generation that is compliant to the Ockr standard, to ensure that the information can be extracted from the document and processed by a machine. The second iteration will also provide the registration process. This could be the more challenging part, as it has a dependency to the Cardano blockchain and needs ADA to submit the transaction.

5.6 Ockr App

The Ockr app is a mobile application that can be used to validate information hashes. It will be available as an open source project under the Ockr umbrella. The app will be available for Android and iOS. Ionic will be used to build the app. The OCR model will be provided using ONNX.js and a model from the Ockr model zoo.

5.7 Chain Submitter Service

The chain submitter service is a backend service that can be used to submit transactions to the Cardano blockchain. It will be available as part of the Ockr api repository. Java Spring Boot will be used to implement the service. The service will also be available as a Docker image.