

# Axiomatic semantics of docxwatermarker

Fabio F.G. Buono

2026

Version 2

Licensed under CC BY 4.0

## 1 Introduction

A `.docx` file is a ZIP archive, and under `word/media/` the archive carries the raw bytes of every image anchored to the document. Replacing one of these images with other bytes and rewriting the archive yields a Word document in which the image appears different and every other piece of the file stays in place. The `docxwatermarker` library does this work, and it serves the developer who prepares a template in Word and wants to produce one personalized copy per recipient, for instance to stamp confidential documents.

These notes give an axiomatic specification of the library. They cover the library's side of a specification in two parts, whose other half, an operational semantics of the command-line tool, is the companion document *Operational semantics of the docxwatermarker command-line tool*. For each public operation they state a Hoare triple [1] characterizing the operation's behaviour, and from the triples they derive the properties that hold across several operations. Under this lens the library becomes a collection of contracts on bytes, archives, images, and ZIP metadata. The notation follows Mancarella's notes [5], with the standard references to Winskel [3] for the foundations and to Meyer [4] for the connection with the executable contracts the code already carries in its `require` and `ensure` primitives.

In the cited texts the triples specify constructs of a programming language, such as assignment or the while loop. Here each triple specifies a function call, and the host language, Python, is treated as an oracle.

## 2 How to read these notes

The Hoare-triple formalism usually applies to programs of some significance, such as critical systems, controllers, synchronization primitives, or, in textbook contexts, to small artificial code fragments. A document-personalization library falls in neither category. The code declares a partial contract through its `require` and `ensure` primitives, and making that contract explicit at the mathematical level closes a loop the code opens only halfway, because every runtime assertion gains a corresponding formal statement, and every property stated here should map to a `require` or `ensure` in the code. Writing the specification also forced us to ask questions about the library's behaviour that the code alone does not push us to ask, and four corrections in version 0.1 came out of this exercise.

The notation  $\{P\}c\{Q\}$  has its standard reading. If the assertion  $P$  holds in the state of the client program's variables before the call to  $c$ , then the assertion  $Q$  holds afterwards, where in  $Q$  the symbol  $y$  denotes the value just returned and uppercase variables such as  $X$  denote the values their lowercase counterparts held on entry. The triples are in total correctness, and

they include termination of  $c$ . Operations that may raise exceptions are treated alongside the main triple, with the conditions that trigger each exception.

Lowercase variables such as  $x, y, t, p$  are client-program variables, and uppercase variables such as  $X, Y, T, P$  are specification variables fixing the entry values so that they can be referred to in the postcondition, following the convention of §6 of [5]. Fixed-width identifiers such as `Template.open` or `replace_image` are the names of classes and methods in the Python code, while mathematical italic identifiers such as *detect*, *img*, *inMedia* are names of mathematical functions that model the library’s internal operations.

### 3 Preliminaries

The sets used in the following are treated in the intuitive notion of collection, outside any formal set theory.

We write **Bytes** for the set of finite octet sequences,  $\text{Bytes} = \{0, 1, \dots, 255\}^*$ . For  $b \in \text{Bytes}$ ,  $|b|$  denotes the length in octets and  $b_1 \cdot b_2$  the concatenation.

We write **Name** for the set of strings used as internal entry names in a ZIP archive, and on **Name** we assume no structure beyond equality and the function *lower* : **Name**  $\rightarrow$  **Name** returning the lowercase form of a name.

We write **Format** for the set of image formats recognized by the library,

$$\text{Format} = \{\text{png}, \text{jpeg}, \text{gif}, \text{bmp}, \text{tiff}, \text{webp}, \perp\}.$$

The value  $\perp$  stands for “unrecognized format”.

#### Entries and archives

An entry in a ZIP archive is a pair of content bytes and a collection of metadata (timestamp, compression method, POSIX permissions, miscellaneous attributes). The internal structure of the metadata plays no role in the triples, and we denote by *Meta* their domain treated as opaque. Entries produced by a modifying operation of the library are distinguished from original entries by a *modification flag* taking values `tt` or `ff`.

**Definition 3.1.** **Entry** is the set of quadruples  $e = (n, d, m, \mu)$  with  $n \in \text{Name}$ ,  $d \in \text{Bytes}$ ,  $m \in \text{Meta}$ ,  $\mu \in \{\text{tt}, \text{ff}\}$ , with components denoted  $e.name$ ,  $e.data$ ,  $e.meta$ ,  $e.mod$ .

**Definition 3.2.** An archive is a partial function  $A : \text{Name} \rightarrow \text{Entry}$  with finite domain, such that for every  $n \in \text{dom}(A)$  we have  $A(n).name = n$ . We write **Archive** for the set of archives.

The order of entries, although relevant in practice (the library takes care to preserve it), never enters into the triples, so we leave it implicit in the phrase “partial function with finite domain”.

On an archive we define the update operation. Given  $A \in \text{Archive}$ ,  $n \in \text{dom}(A)$  and  $e \in \text{Entry}$  with  $e.name = n$ , we write  $A[e/n]$  for the archive defined by

$$A[e/n](n') = \begin{cases} e & \text{if } n' = n; \\ A(n') & \text{otherwise,} \end{cases}$$

in direct analogy with Mancarella’s function-update operator [5, §2].

#### Images

An image is the logical view of an entry whose content is recognized as one of the formats in  $\text{Format} \setminus \{\perp\}$ . Beyond name and format, the data of interest are width, height, size in bytes, and a flag indicating whether the image carries the *marker* produced by `make_marker_png`, which is the mechanism through which the `by_marker` matcher identifies the right image in a template.

**Definition 3.3.** *Image* is the set of sextuples  $i = (n, f, w, h, s, k)$  with  $n \in \mathbf{Name}$ ,  $f \in \mathbf{Format} \setminus \{\perp\}$ ,  $w, h, s \in \mathbb{N}$ ,  $k \in \{\text{tt}, \text{ff}\}$ , subject to the constraint

$$k = \text{tt} \implies f = \text{png}, \quad (1)$$

with components accessible as  $i.name, i.f, i.w, i.h, i.s, i.mark$ .

Constraint (1) reflects the fact that the marker mechanism is based on textual `tEXt` chunks, a feature available only in the PNG format.

The total functions used in the following have signatures

$$\begin{aligned} detect &: \mathbf{Bytes} \rightarrow \mathbf{Format}, \\ dim &: \mathbf{Bytes} \rightarrow \mathbb{N} \times \mathbb{N}, \\ isMarker &: \mathbf{Bytes} \rightarrow \{\text{tt}, \text{ff}\}, \\ inMedia &: \mathbf{Name} \rightarrow \{\text{tt}, \text{ff}\}. \end{aligned}$$

The *detect* function identifies the format of a byte sequence from its leading *magic bytes* and returns  $\perp$  when no signature matches. The *dim* function returns pixel dimensions, or  $(0, 0)$  when the bytes are not decodable as an image. The *isMarker* predicate returns `tt` iff the first 512 bytes begin with the PNG signature and contain the marker tag. The *inMedia* predicate holds iff *lower*( $n$ ) starts with the prefix "word/media/".

Given an entry  $e$  with  $detect(e.data) \neq \perp$ , we set

$$img(e) = (e.name, detect(e.data), dim(e.data)_1, dim(e.data)_2, |e.data|, isMarker(e.data)).$$

The well-definedness of *img* with respect to constraint (1) is immediate, since *isMarker* returns `tt` only for byte sequences starting with the PNG signature, on which *detect* returns `png`.

## Templates

A template is an archive, and we set  $\mathbf{Template} = \mathbf{Archive}$ . In practical use a template is an archive carrying a number of mandatory entries (`[Content_Types].xml`, `word/document.xml`, and so on), but these requirements do not enter into the library's triples.

On a template we define the function  $images : \mathbf{Template} \rightarrow \mathbf{Image}^*$  enumerating the images through

$$images(t) = [img(t(n)) \mid n \in dom(t), inMedia(n), detect(t(n).data) \neq \perp],$$

where the notation  $[\cdot \mid \cdot]$  denotes a sequence built by filtering (the concrete order is the iteration order on the archive's domain, and none of our properties depend on this order).

The matchers, exposed by the library as strategies for selecting an image inside a template, are partial functions  $\mu : \mathbf{Image}^* \rightarrow \mathbf{Name}$ . The three concrete matchers  $\mu_{\text{name}}^p, \mu_{\text{mark}}, \mu_{\text{auto}}$  will be defined in §5.5, where they serve as parameters to the triples for `replace_image`.

## 4 States and assertions

The client program that uses the library is written in Python and manipulates ordinary variables. To reason about the behaviour of the library's operations we abstract from the client language and consider only the effect on the state of its variables.

Let  $\mathbf{Var}$  be a set of Python variable identifiers. We write

$$\mathbf{Val} = \mathbf{Template} \cup \mathbf{Image}^* \cup \mathbf{Bytes} \cup \mathbf{Name} \cup \mathbb{Z} \cup \{\text{tt}, \text{ff}\}$$

for the set of values our variables can take in these arguments. A *state* is a function

$$\sigma : \mathbf{Var} \rightarrow \mathbf{Val}$$

that maps each assigned variable to its value, and stays undefined on variables not yet assigned. We write **State** for the set of all states.

State update is defined by

$$\sigma[v/x](y) = \begin{cases} v & \text{if } y = x, \\ \sigma(y) & \text{otherwise,} \end{cases}$$

as the function-update operator of Paragraph 3.

An *assertion* is a first-order predicate-calculus formula over the variables in **Var**. We use letters  $P, Q, R, \dots$  for assertions. The constants in the assertions are those of the domains introduced in Paragraph 3, and terms like  $t.images$  (with  $t \in \mathbf{Var}$ ) are state-dependent expressions. We extend the usual infix notation to the objects of our domain, so  $\in$  stands for sequence membership,  $=$  for equality,  $|s|$  for the length of a sequence.

We write  $\sigma \models P$  for “ $P$  is true in state  $\sigma$ ”, in the standard inductive definition on the structure of  $P$ . We also use specification variables (uppercase letters  $X, Y, A, B, \dots$ ) to fix generic values on which to parametrize preconditions and postconditions, following the convention of §6 of [5].

## Hoare triples for library operations

A *Hoare triple*  $\{P\}c\{Q\}$  is satisfied if, for every state  $\sigma$  such that  $\sigma \models P$ , executing  $c$  from  $\sigma$  terminates in a state  $\sigma'$  with  $\sigma' \models Q$ . In our setting,  $c$  is always a library function call of the form

$$y \leftarrow op(\bar{a}),$$

where  $op$  is an operation,  $\bar{a}$  is the tuple of arguments, and  $y \in \mathbf{Var}$  receives the returned value, when the operation produces one. In the postcondition  $Q$  the symbol  $y$  denotes the value just returned. The variables denoting the arguments may be annotated through the *specification variables* convention, by writing  $X$  in  $Q$  to refer to the value the variable  $x$  held in  $\sigma$  before the call.

A triple may be weakened or strengthened through the classical rule

$$\frac{P \Rightarrow P' \quad \{P'\}c\{Q'\} \quad Q' \Rightarrow Q}{\{P\}c\{Q\}} \quad (\text{pre-post})$$

which we use as an auxiliary tool in the proofs of Paragraph 6.

**Example 4.1.** Consider the triple

$$\begin{aligned} & \{ p \in \mathbf{Bytes} \wedge p \text{ is a valid ZIP} \} \\ & t \leftarrow \mathbf{Template.from\_bytes}(p) \\ & \{ t \in \mathbf{Template} \wedge toBytes(t) = p \}. \end{aligned}$$

The precondition states that  $p$  is a byte sequence forming a valid ZIP archive, where “ $p$  is a valid ZIP” is the predicate realized by the Python `zipfile` library (we do not define it here). The postcondition asserts that, after the call, the variable  $t$  holds a template and that the serialization of that template via `toBytes` returns the same input bytes. The fact that the Python serializer preserves the bytes of a ZIP archive read and written back without changes (modulo details discussed in the paragraph on idempotence) is the reason why the postcondition carries the conjunct  $toBytes(t) = p$  in addition to the typing  $t \in \mathbf{Template}$ .

## 5 Operation specifications

The triples for the library’s operations share a convention on exceptions, which the library shapes around the hierarchy

$$\text{DocxWatermarkerError} \supseteq \{ \text{ImageNotFoundError}, \text{MultipleImagesError}, \\ \text{FormatMismatchError}, \text{PDFConversionError}, \text{InvariantError} \},$$

defined in the `errors` module. The syntax we use to indicate that the call  $c$ , executed in a state satisfying  $P$ , may terminate in a state satisfying  $Q$  or raise one of the exceptions  $E_i$  is

$$\{P\} c \{Q\} \mid E_1 \text{ if } R_1, \dots, E_n \text{ if } R_n,$$

with the convention that  $E_i$  is raised on the first matching condition  $R_i$ .

### 5.1 Template constructors

The operation `Template.from_bytes(p)` builds a template from raw bytes. The operation `Template.open(l)` does the same from a filesystem path.

**Template.from\_bytes.**

$$\{ p \in \text{Bytes} \wedge \text{isValidZip}(p) \} \\ t \leftarrow \text{Template.from\_bytes}(p) \\ \{ t \in \text{Template} \wedge \text{toBytes}(t) \approx_{\text{sem}} p \}$$

The predicate  $\text{isValidZip} : \text{Bytes} \rightarrow \{\text{tt}, \text{ff}\}$  is the one realized by `zipfile.ZipFile`, and the relation  $\approx_{\text{sem}}$  is the *semantic equivalence* between bytes of ZIP archives, defined by  $b_1 \approx_{\text{sem}} b_2$  iff the two serialized archives have the same set of names, the same per-entry contents and the same per-entry metadata, modulo possible differences in the Local File Header padding bytes computed by the serializer. The distinction between *semantic* equivalence ( $\approx_{\text{sem}}$ , which we use throughout) and *physical* equivalence (byte-for-byte identity, which holds only with the `reproducible=True` flag) is discussed in the project README under the heading *byte-perfectness*. When  $p$  is not a valid ZIP, the behaviour is not defined by the specification, and the code raises `BadZipFile`, a Python standard library exception that lives outside the `DocxWatermarkerError` hierarchy.

**Template.open.** Opening a template from the filesystem is the composition of file reading and `from_bytes`, with the addition that a file-reading failure is tracked as `FileNotFoundError`.

$$\{ \ell \in \text{Path} \wedge \text{fileAt}(\ell) = p \wedge \text{isValidZip}(p) \} \\ t \leftarrow \text{Template.open}(\ell) \\ \{ t \in \text{Template} \wedge \text{toBytes}(t) \approx_{\text{sem}} p \} \mid \text{FileNotFoundError} \text{ if } \neg \text{exists}(\ell)$$

$\text{fileAt}(\ell)$  denotes the bytes of the file at the system address  $\ell$ , and  $\text{exists}(\ell)$  the file-existence predicate.

### 5.2 Template inspection

**Template.list\_images.** Returns the sequence of images contained in a template, according to the function *images* of Paragraph 3.

$$\{ t \in \text{Template} \} \\ \bar{i} \leftarrow t.\text{list\_images}() \\ \{ \bar{i} = \text{images}(T) \wedge t = T \}$$

The conjunct  $t = T$  in the postcondition expresses the immutability of the template, and will return in Paragraph 6 as a proved property.

### 5.3 Image replacement

The operation  $t.\text{replace\_image}(x, \mu, v)$  is the library's pivot. It replaces, inside the template  $t$ , the image selected by the matcher  $\mu$  with the bytes obtained from a normalization of the input  $x$ , and it returns a new template. The parameter  $v \in \{\text{tt}, \text{ff}\}$  is the *validate* flag activating image decode validation.

The triple uses two auxiliary functions. The normalization  $normalize : \text{Val} \rightarrow \text{Bytes}$  is defined on the union of **Bytes**, *Path* and *PILImage* (the set of Pillow images), and returns the input bytes when they are already in **Bytes**, the file bytes for a *Path*, the bytes of a PNG encoding for a Pillow image. The decoder  $canDecode : \text{Bytes} \rightarrow \{\text{tt}, \text{ff}\}$  returns **tt** when Pillow opens and verifies the bytes as a valid image, **ff** otherwise.

**Template.replace\_image.**

$$\left\{ \begin{array}{l} t \in \text{Template} \wedge t = T \wedge \mu \in \text{Matcher} \wedge x \in \text{dom}(\text{normalize}) \\ \wedge \mu(\text{images}(t)) = n \wedge \text{detect}(\text{normalize}(x)) = T(n).f \\ \wedge (v = \text{tt} \Rightarrow \text{canDecode}(\text{normalize}(x))) \end{array} \right\}$$

$$t' \leftarrow t.\text{replace\_image}(x, \mu, v)$$

$$\left\{ \begin{array}{l} t' \in \text{Template} \wedge t' = T[e/n] \wedge t = T \\ \text{where } e.\text{data} = \text{normalize}(x), e.\text{name} = n, e.\text{meta} = T(n).\text{meta}, e.\text{mod} = \text{tt} \end{array} \right\}$$

The call raises

- **ImageNotFoundError** if  $\mu(\text{images}(t))$  is undefined because of no candidates;
- **MultipleImagesError** if  $\mu(\text{images}(t))$  is undefined because of ambiguity (more than one candidate);
- **FormatMismatchError** if  $\text{detect}(\text{normalize}(x)) \neq T(n).f$ , namely the new input's format differs from the target's;
- **FormatMismatchError** if  $v = \text{tt}$  and  $\neg \text{canDecode}(\text{normalize}(x))$ , namely the declared format is correct but the bytes are not decodable.

The precondition asks for  $t = T$ , with  $T$  a specification variable, because the postcondition itself asserts  $t = T$  and needs a name for the entry value of  $t$ . The returned template is  $T[e/n]$ , the update of  $T$  on the single name  $n$  with the new entry  $e$ , and the other entries are left untouched (explicit *frame*). The entry  $e$  inherits the metadata  $T(n).\text{meta}$ , comprising compression method, permissions and other attributes, with  $e.\text{data}$  and the flag  $e.\text{mod} = \text{tt}$  as the only changes. The clause  $v = \text{tt} \Rightarrow \text{canDecode}(\text{normalize}(x))$  is a precondition that activates only when the caller asks for validation, reflecting the code's choice of an opt-in validation with default  $v = \text{tt}$ .

### 5.4 Serialization

The operations  $t.\text{to\_bytes}()$  and  $t.\text{save}(\ell)$  produce the bytes of a template, respectively in memory and on disk.

`Template.to_bytes`.

$$\begin{aligned} & \{ t \in \text{Template} \wedge t = T \} \\ & \quad b \leftarrow t.\text{to\_bytes}() \\ & \{ b \in \text{Bytes} \wedge \text{isValidZip}(b) \wedge \text{fromBytes}(b) = T \wedge t = T \} \end{aligned}$$

where *fromBytes* is the abstract function that inverts *toBytes* up to semantic equivalence. Reapplying the reading to the produced bytes returns *T*, and the round-trip roundness *toBytes*  $\circ$  *fromBytes*  $\approx$  id underlies the *byte-perfectness* promise that the library makes in the README.

`Template.save`.

$$\begin{aligned} & \{ t \in \text{Template} \wedge t = T \wedge \ell \in \text{Path} \wedge \text{writable}(\ell) \} \\ & \quad \ell' \leftarrow t.\text{save}(\ell) \\ & \{ \ell' \in \text{Path} \wedge \text{fileAt}(\ell') = t.\text{to\_bytes}() \wedge t = T \} \mid \text{OSError if } \neg \text{writable}(\ell) \end{aligned}$$

*writable*( $\ell$ ) is the writability predicate on the filesystem at the address  $\ell$ .

## 5.5 Matchers

The constructors of `ImageMatcher` are pure functions, and what characterizes each one is the partial function  $\mu : \text{Image}^* \rightarrow \text{Name}$  that the produced matcher computes.

`ImageMatcher.by_filename`( $p$ ).

$$\begin{aligned} & \{ p \in \text{Name} \} \\ & \quad \mu \leftarrow \text{ImageMatcher.by\_filename}(p) \\ & \{ \mu \in \text{Matcher} \wedge \forall \bar{i} \in \text{Image}^* . \mu(\bar{i}) = \begin{cases} i.\text{name} & \text{if } \exists! i \in \bar{i} . i.\text{name} = p, \\ \text{undef.} & \text{otherwise} \end{cases} \} \end{aligned}$$

`ImageMatcher.by_marker`( $\cdot$ ).

$$\begin{aligned} & \{ \text{true} \} \\ & \quad \mu \leftarrow \text{ImageMatcher.by\_marker}() \\ & \{ \mu \in \text{Matcher} \wedge \forall \bar{i} \in \text{Image}^* . \mu(\bar{i}) = \begin{cases} i.\text{name} & \text{if } \exists! i \in \bar{i} . i.\text{mark} = \text{tt}, \\ \text{undef.} & \text{otherwise} \end{cases} \} \end{aligned}$$

`ImageMatcher.auto`( $s$ ). A monoidal combinator of the previous two, with a second heuristic stage parametric in  $s \in \mathbb{N}$ , the minimum side of the candidate square.

$$\begin{aligned} & \{ s \in \mathbb{N} \} \\ & \quad \mu \leftarrow \text{ImageMatcher.auto}(s) \\ & \{ \mu \in \text{Matcher} \wedge \forall \bar{i} . \mu(\bar{i}) = \mu_{\text{mark}}(\bar{i}) \text{ or else } \mu_{\text{heur}}^s(\bar{i}) \} \end{aligned}$$

where *or else* is the standard combinator “try the first, if undefined try the second”, and

$$\mu_{\text{heur}}^s(\bar{i}) = \begin{cases} i.\text{name} & \text{if } \exists! i \in \bar{i} . i.f = \text{png} \wedge i.w = i.h \wedge i.w \geq s, \\ \text{undef.} & \text{otherwise.} \end{cases}$$

## 5.6 Image generation

The function `make_marker_png` generates a PNG placeholder carrying the marker tag. The function `make_text_watermark` generates a text watermark.

`make_marker_png(s)`.

$$\begin{aligned} & \{ s \in \mathbb{N} \wedge s \geq 1 \} \\ & d \leftarrow \text{make\_marker\_png}(s) \\ & \{ d \in \text{Bytes} \wedge \text{detect}(d) = \text{png} \wedge \text{isMarker}(d) = \text{tt} \wedge \text{dim}(d) = (s, s) \} \end{aligned}$$

`make_text_watermark( $\bar{\lambda}$ ,  $s$ ,  $\rho$ ,  $\chi$ , ...)`. The function produces a square PNG of side  $s$  with the text lines  $\bar{\lambda}$ , rotated by  $\rho$  degrees, in colour  $\chi$ . The dictionary of presets, denoted *PRESETS*, is defined in the code.

$$\begin{aligned} & \left\{ \begin{array}{l} (\bar{\lambda} \in \text{String}^* \wedge \bar{\lambda} \neq \varepsilon \wedge \exists \lambda \in \bar{\lambda}. \text{trim}(\lambda) \neq \varepsilon) \vee \\ (\text{preset} \in \text{dom}(\text{PRESETS})) \wedge s \in \mathbb{N} \wedge s \geq 1 \end{array} \right\} \\ & d \leftarrow \text{make\_text\_watermark}(\bar{\lambda}, s, \rho, \chi, \dots) \\ & \{ d \in \text{Bytes} \wedge \text{detect}(d) = \text{png} \wedge \text{dim}(d) = (s, s) \} \\ & | \text{ValueError} \text{ if both conjuncts of the precondition fail} \end{aligned}$$

The specification does not characterize the pictorial content of the produced image, since doing so would require a model of the rendering process. The relevant property is that  $d$  is a valid PNG of size  $(s, s)$ .

## 5.7 PDF conversion

`find_libreoffice()`. A pure function returning the path to `soffice`, or  $\perp$  when no path is found.

$$\begin{aligned} & \{ true \} \\ & r \leftarrow \text{find\_libreoffice}() \\ & \{ (r \in \text{Path} \wedge \text{executable}(r)) \vee r = \perp \} \end{aligned}$$

`to_pdf( $\ell_{src}$ ,  $\ell_{dst}$ ,  $\tau$ ,  $\beta$ )`. DOCX-to-PDF conversion goes through LibreOffice in headless mode. The parameters are the source path  $\ell_{src}$ , the destination path  $\ell_{dst}$  (possibly  $\perp$  for a default derived from the source name), the timeout  $\tau$ , and the binary override  $\beta$ .

$$\begin{aligned} & \left\{ \begin{array}{l} \ell_{src} \in \text{Path} \wedge \text{exists}(\ell_{src}) \wedge \tau > 0 \\ \wedge (\beta \neq \perp \Rightarrow \text{executable}(\beta)) \wedge (\beta = \perp \Rightarrow \text{find\_libreoffice}() \neq \perp) \end{array} \right\} \\ & \ell' \leftarrow \text{to\_pdf}(\ell_{src}, \ell_{dst}, \tau, \beta) \\ & \{ \ell' \in \text{Path} \wedge \text{exists}(\ell') \wedge \text{isPDF}(\text{fileAt}(\ell')) \} \end{aligned}$$

The exceptions the function may raise are all subtypes of `PDFConversionError`, parametric in a *reason* field identifying the cause.

- `PDFConversionError(reason = not_found)` if the precondition on  $\beta$  or `find_libreoffice` fails;
- `PDFConversionError(reason = cannot_create_outdir)` if the directory containing  $\ell_{dst}$  cannot be created;

- `PDFConversionError(reason = timeout)` if the conversion exceeds  $\tau$  seconds;
- `PDFConversionError(reason = nonzero_exit)` if `soffice` exits with a non-zero status;
- `PDFConversionError(reason = no_output)` if `soffice` exits cleanly but the PDF file is not produced.

The postcondition `isPDF(fileAt( $\ell'$ ))` is the only one in the document that does not descend from the library's code, since it depends on the correct functioning of LibreOffice, which we treat as an oracle. The code does ensure that every subprocess failure is tracked as one of the variants of `PDFConversionError` above.

## 6 System invariants

The triples of the previous paragraph speak of one operation at a time. Some properties, however, cross several operations. Template immutability is one of these, and it holds for every operation taking a template on entry without building one from bytes.

**Property 6.1.** *For every  $t \in \text{Template}$  and every operation  $o$  in*

$$\{\text{list\_images}, \text{replace\_image}, \text{to\_bytes}, \text{save}\},$$

*after the execution of  $o$  on  $t$  the value of  $t$  is the one it held on entry.*

*Proof.* The four triples in Paragraph 5 all contain  $t = T$  in the postcondition, with  $T$  denoting the entry value of  $t$ .  $\square$

In the code the property is guaranteed by the type, because the constructor of `Template` freezes the dictionary with `types.MappingProxyType` (in `Template.__init__`), and `replace_image` builds a new dictionary from the original without mutating it. A related invariant fixes the template's domain under `replace_image`.

**Property 6.2.** *For every  $t \in \text{Template}$  and every legitimate call to `replace_image`, the returned template  $t'$  satisfies  $\text{dom}(t') = \text{dom}(t)$ .*

*Proof.* From the postcondition of `replace_image`,  $t' = T[e/n]$  with  $n \in \text{dom}(T)$ . By the definition of archive update,  $\text{dom}(T[e/n]) = \text{dom}(T)$ . From  $t = T$  the thesis follows.  $\square$

An `ensure` (spec=I2) in `replace_image` checks the equality of the two namesets and raises `InvariantError` when the property is violated. It is the only structural `ensure` in the library, because adding or removing an entry from an OOXML archive would break the document. A similar guarantee holds on the format of the images, since the file `[Content_Types].xml` declares the MIME type of each entry, and replacing a PNG with a JPEG without updating that file would yield an unreadable document.

**Property 6.3.** *For every  $t \in \text{Template}$  and every legitimate call to `replace_image(x,  $\mu$ ,  $v$ )` with  $\mu(\text{images}(t)) = n$ , the returned template  $t'$  satisfies  $\text{img}(t'(n)).f = \text{img}(t(n)).f$ .*

*Proof.* The precondition contains  $\text{detect}(\text{normalize}(x)) = T(n).f$ . The entry  $e$  in the postcondition has  $e.data = \text{normalize}(x)$ , so  $\text{detect}(e.data) = T(n).f$ . From  $t'(n) = e$  follows  $\text{img}(t'(n)).f = \text{img}(t(n)).f$ .  $\square$

The constraint is in `replace_image`, before the substitution takes place, with `FormatMismatchError` when the precondition fails. The matcher `by_marker` has an unambiguous semantics only when `make_marker_png` has been called once on the template of interest.

**Property 6.4.** *If in  $\text{images}(t)$  there is at most one image  $i$  with  $i.mark = \text{tt}$ , then  $\mu_{\text{mark}}$  is defined on  $\text{images}(t)$  iff such an image exists, and it returns its name.*

*Proof.* The definition of  $\mu_{\text{mark}}(\bar{i})$  yields  $i.name$  when a unique  $i \in \bar{i}$  exists with  $i.mark = tt$ , and remains undefined otherwise. The uniqueness hypothesis removes the “more than one” case.  $\square$

The library refrains from checking the convention, because inspecting the PNG metadata of every image of every template would have a hard-to-justify cost, and any violation surfaces as a `MultipleImagesError` raised by `replace_image`. Reading and writing of a template form an inverse pair, modulo the padding of the ZIP Local File Header.

**Property 6.5.** *For every  $t \in \text{Template}$ ,  $\text{fromBytes}(\text{toBytes}(t)) \approx_{\text{sem}} t$ .*

*Proof.* From the postcondition of `Template.to_bytes`,  $\text{fromBytes}(b) = T$  for the returned bytes  $b$ . From the postcondition of `Template.from_bytes`,  $\text{toBytes}(t') \approx_{\text{sem}} p$  for every  $t'$  obtained from bytes  $p$ . With  $b = \text{toBytes}(T)$  and  $t' = \text{fromBytes}(b)$  we get  $\text{toBytes}(t') \approx_{\text{sem}} \text{toBytes}(T)$ , hence  $t' \approx_{\text{sem}} T$ .  $\square$

A runtime comparison would be too expensive and the property is covered by tests on a sample of real templates in `tests/test_zipops.py`. Under the `reproducible=True` flag of `_zipops.write_zip` the relation  $\approx_{\text{sem}}$  collapses to byte equality, because the timestamps are reset to the DOS minimum, and the property becomes a literal equality.

The triples of Paragraph 5 are in total correctness, so they include termination. The large majority of the library’s functions contain no loops, and they terminate by the termination of the system and Pillow calls. The only non-trivial loop is the binary search for the largest  $s$  in  $[s_{\text{min}}, s_{\text{max}}]$  such that the text rendered at font size  $s$  fits inside the box, in `_find_max_font_size` in `watermark.py`. The termination function is  $s_{\text{max}} - s_{\text{min}}$ , which halves at each iteration, and the proof follows Mancarella’s while rule [5, §4.7].

## 7 Mapping to the code

The operation specifications do not refer to the code, and the invariants of Paragraph 6 only in passing. Each clause is tied to the code by a `spec=` identifier carried on the realizing `require` or `ensure`, not by a line number. Line numbers rot as the code shifts, whereas the identifier moves with the check. The test `test_spec_crossref.py` verifies that the identifiers used in the code, the catalogue in `_invariants.SPEC_REALIZATIONS`, and the tables below stay in agreement. The Location column names the file and the Python identifier. The Assertions column gives the `spec=` where one is present.

Formal operation	Location	Assertions
<code>Template.open</code>	<code>core.py</code>	—
<code>Template.from_bytes</code>	<code>core.py</code>	<code>require spec=P:from_bytes</code>
<code>Template.list_images</code>	<code>core.py</code>	—
<code>Template.replace_image</code>	<code>core.py</code>	<code>ensure spec=I2</code>
<code>Template.save</code>	<code>core.py</code>	—
<code>Template.to_bytes</code>	<code>core.py</code>	—
$\mu_{\text{fname}}$	<code>core.py</code> ( <code>by_filename</code> )	<code>require spec=P:by_filename</code>
$\mu_{\text{mark}}$	<code>core.py</code> ( <code>by_marker</code> )	—
$\mu_{\text{auto}}$	<code>core.py</code> ( <code>auto</code> )	—
<code>normalize</code>	<code>core.py</code> ( <code>_normalize_image_input</code> )	—
<code>detect</code>	<code>_imagedetect.py</code> ( <code>_detect_format</code> )	—
<code>images</code>	<code>_imagedetect.py</code> ( <code>list_images</code> )	—
<code>dim</code>	<code>_imagedetect.py</code> ( <code>_read_dimensions</code> )	—
<code>make_marker_png</code>	<code>_imagedetect.py</code>	—
<code>isMarker</code>	<code>_imagedetect.py</code> ( <code>is_marker_image</code> )	—
<code>make_text_watermark</code>	<code>watermark.py</code>	—
<code>find_libreoffice</code>	<code>pdf.py</code>	—
<code>to_pdf</code>	<code>pdf.py</code>	—
<code>toBytes</code> (internal)	<code>_zipops.py</code> ( <code>write_zip</code> )	—
<code>fromBytes</code> (internal)	<code>_zipops.py</code> ( <code>read_zip_entries</code> )	—

<b>Invariant</b>	<b>Realization</b>	<b>Enforced by</b>
I1 (immutability)	MappingProxyType in <code>Template.__init__</code>	structure
I2 (nameset preservation)	<code>ensure(spec="I2")</code> in <code>replace_image</code>	ensure
I3 (format preservation)	preventive <code>FormatMismatchError</code>	structure
I4 (marker uniqueness)	<code>by_marker</code> uniqueness convention	convention
I5 (semantic round-trip)	property tests in <code>test_zipops.py</code>	tests

The table above reflects version 0.1. Version 0.2 will bring changes, in particular the exposure of the `reproducible` flag through `Template.save`, and this table will be updated.

## References

- [1] C. A. R. Hoare. An axiomatic basis for computer programming. *Communications of the ACM*, 12(10):576–580, 1969.
- [2] E. W. Dijkstra. Guarded commands, nondeterminacy and formal derivation of programs. *Communications of the ACM*, 18(8):453–457, 1975.
- [3] G. Winskel. *The Formal Semantics of Programming Languages: An Introduction*. MIT Press, Cambridge, MA, 1993. ISBN 0-262-73103-7. In particular Chapter 6, *Axiomatic semantics*.
- [4] B. Meyer. Applying “design by contract”. *IEEE Computer*, 25(10):40–51, 1992.
- [5] P. Mancarella. Note di semantica assiomatica. Lecture notes for the course in Logica per la Programmazione, B.Sc. in Computer Science, University of Pisa, 2008/2009. <https://pages.di.unipi.it/corradini/Didattica/LPP-17/Logica%20di%20Hoare.pdf>