

Paperwork manual

Contents

1	Scanning	4
1.1	Default source	4
1.2	Scanner calibration	4
1.3	Scan resolution	5
1.4	Single scan	6
1.5	From feeder	6
1.6	OCR languages	7
1.6.1	Windows	7
1.6.2	Debian	8
1.6.3	Fedora	8
1.6.4	Ubuntu	8
1.7	OCR enabled / disabled	8
2	Importing	9
2.1	Images	9
2.2	PDF	10
2.3	Many PDFs in one shot	10
3	Labels	11
3.1	Creating new labels	11
3.2	Setting labels on documents	11
3.3	Modifying a label	12
3.4	Deleting a label	12
4	Searching	13
4.1	Simple search	13
4.2	Advanced search	13

5	Viewing	13
5.1	View pages as grid	13
5.2	View pages as list	14
5.3	Zoom level	14
6	Exporting	14
6.1	Document	15
6.2	Page	15
7	Printing	15
8	Copying text	15
9	Editing pages	15
10	Moving pages	16
10.1	inside a document	16
10.2	from a document to another	16
11	Switching to another work directory	16
12	Backup	18
13	Synchronisation	18
13.1	USB key / USB drive	18
13.2	File Synchronization applications	18
13.2.1	DropBox	18
13.2.2	Shared folder	19
14	Encryption	25
14.1	Windows	25
14.2	GNU/Linux	25
14.2.1	Ecryptfs	25
14.2.2	Encfs	26
15	Advanced use and information	26
15.1	Redo OCR	26
15.1.1	On all the documents	26
15.1.2	On one document	26
15.2	Highlight all words	27

15.3	Keyboard shortcuts	27
15.4	Paperwork's files locations	27
15.5	Work directory layout	27
15.5.1	Global organisation	28
15.5.2	hOCR files	29
15.5.3	Label files	29
15.6	Statistics	29
16	Getting support / reporting issues	29
16.1	Diagnostic dialog	29
16.2	Github issue tracker	29
16.3	Mailing-list	29
17	Uninstalling	29
17.1	Windows	30
17.2	GNU/Linux	30

1 Scanning

1.1 Default source

Some scanners have many sources/input. Basic scanners have only one source : Flatbed. Some others have a feeder, allowing them to scan many pages at one.

The settings here is the source to use for single scan.

If you select "Flatbed" here and use the multi-scan dialog (see 1.5 on page 6), Paperwork will automatically switch to the feeder (if one is found ; otherwise the default source is used too).

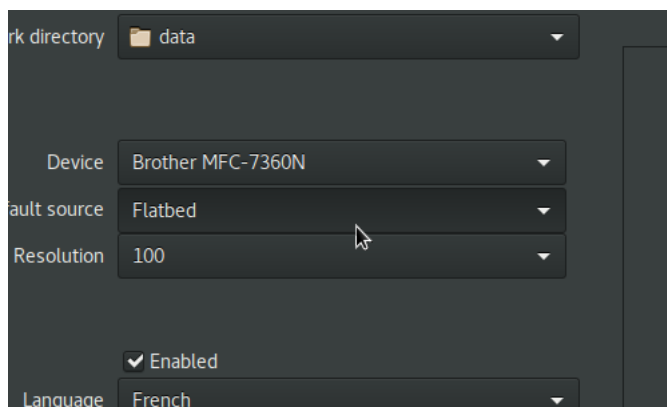


Figure 1: Selecting the default source for scanning



Flatbed only	
Flatbed and feeder	

Table 1: Flatbeds and feeders

1.2 Scanner calibration

Scanners tend to provide images actually bigger than the scanned pages. Since most of the time, you will always scan pages having the same size (A4/Letter usually), Paperwork provides an option called

"scanner calibration". Scanner calibration in Paperwork is simply a pre-cropping of the images coming from the scanner.

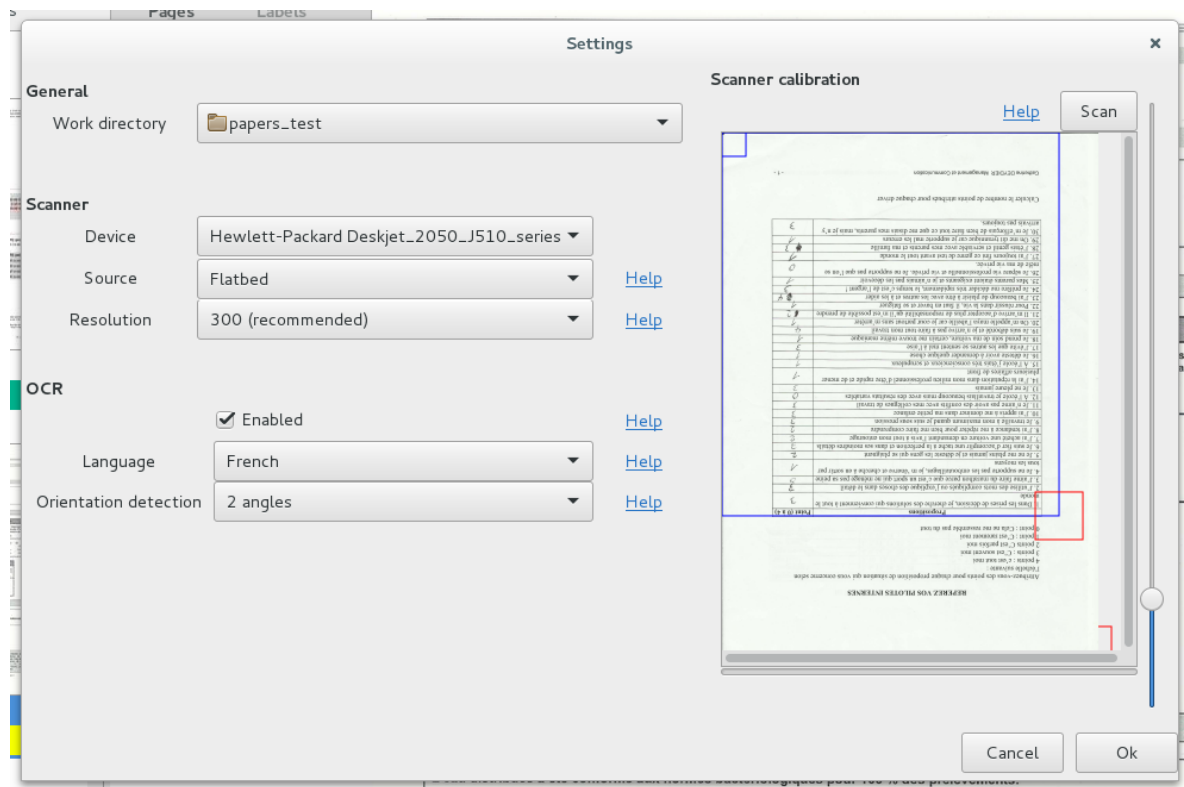


Figure 2: Scanner calibration

1.3 Scan resolution

Scanner resolution defines how detailed the images coming from your scanner must be.

Higher resolutions mean

- longer scans,
- longer OCR,
- more time to display,
- more space used on disk,
- but also better OCR.

Lower resolutions mean

- shorter scans,
- shorter OCR,
- less time to display,

- less space used on disk,
- but also inferior OCR,
- and possibly unreadable image (even by a human).

300 dpi is considered a good trade-off. You may want to reduce it to 200 dpi on slow computers.

1.4 Single scan

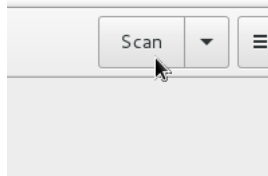


Figure 3: Scanning a single page

Pages are appended to the current document.

1.5 From feeder

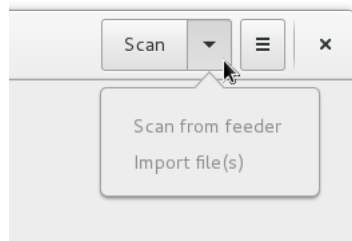


Figure 4: Opening the dialog to scan from the feeder

The option "scan from feeder" is enabled only if Paperwork has detected a feeder on your scanner.

You have to tell Paperwork how many pages go in each document. If you just want Paperwork to scan pages until none are left, you can just specify a huge number of pages (99 for example).

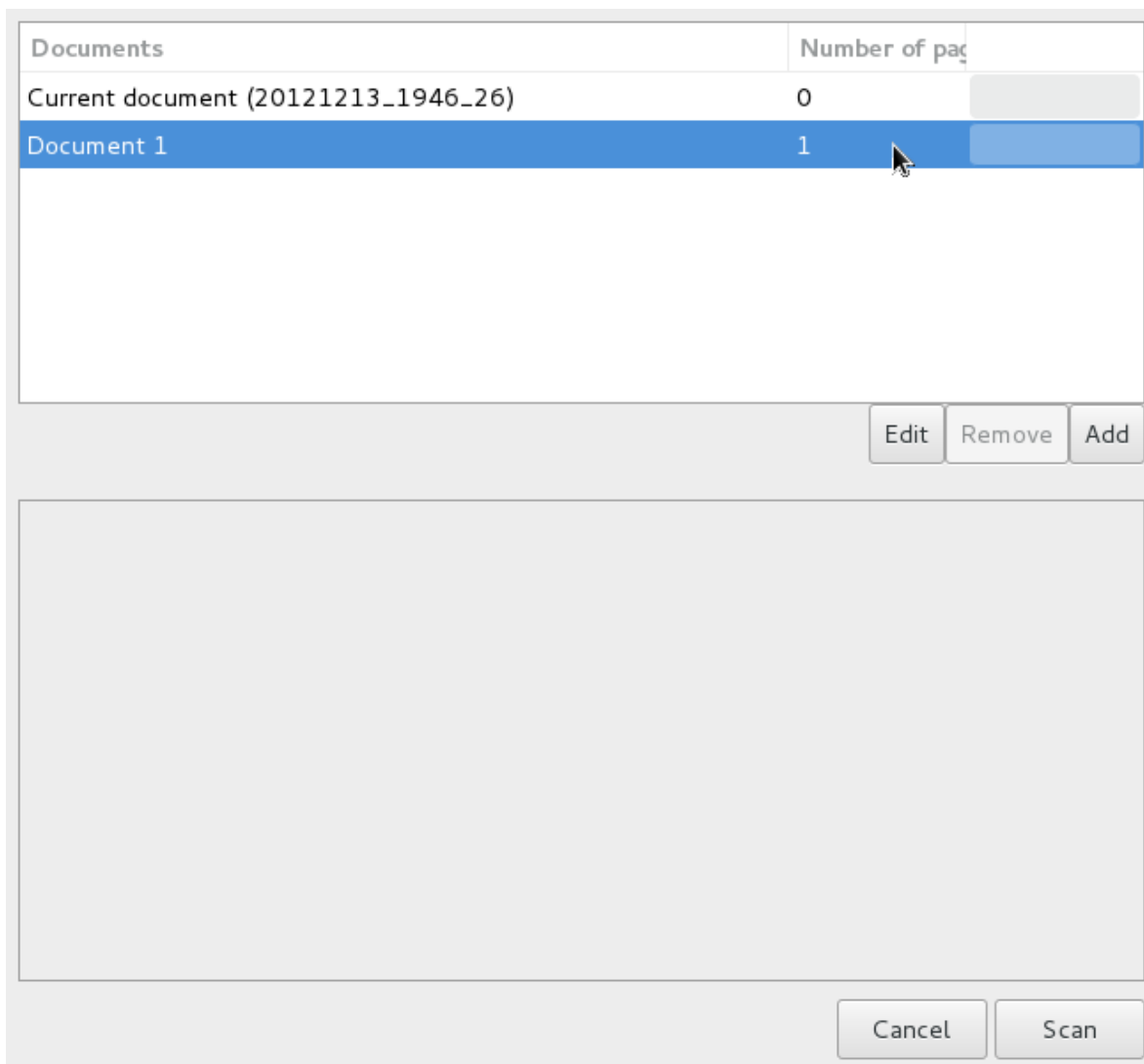


Figure 5: Scanning multiple documents from the feeder

1.6 OCR languages

By default, Paperwork uses Tesseract for the OCR. If unavailable, it falls back on Cuneiform. On Windows, Tesseract is provided with Paperwork.

To get better results, OCR tool need to know the language used in the document(s).

The language available in the settings dialog of Paperwork are those understood by the OCR tool. If your language is not in the list, it means the OCR tool doesn't have the data required to read your language and you must install them.

1.6.1 Windows

Tesseract and all its data files are provided by Paperwork's installer. If a language is not available in the installer, it either means it hasn't been packaged (in which case you can request it), or there is no

data file available yet for this language.

1.6.2 Debian

```
# OCR (Tesseract)
$ sudo apt-get install tesseract-ocr tesseract-ocr-<lang>
```

1.6.3 Fedora

```
# OCR (Tesseract)
$ sudo yum install tesseract tesseract-langpack-<lang>
```

1.6.4 Ubuntu

```
# OCR (Tesseract)
$ sudo apt-get install tesseract-ocr tesseract-ocr-<lang>
```

1.7 OCR enabled / disabled

When you scan a page using Paperwork, Paperwork will immediately run the OCR on it. This process may take a while for each page.

In case you want to scan a lot of pages quickly (for instance, the first time you use Paperwork), OCR can be temporarily disabled.

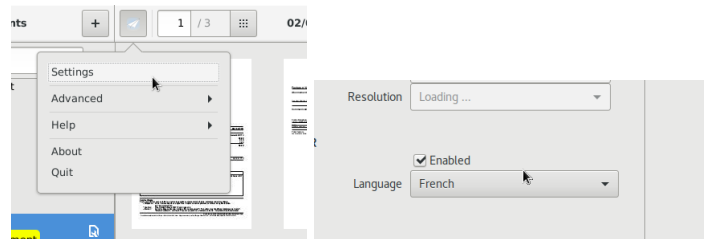


Figure 6: Disabling the OCR

OCR can then be run on all the documents managed by Paperwork in one shot.

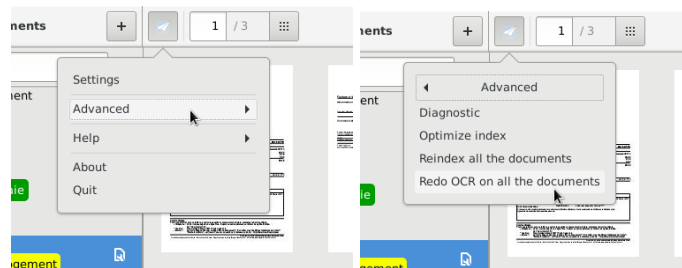


Figure 7: Redo OCR on all the documents

2 Importing

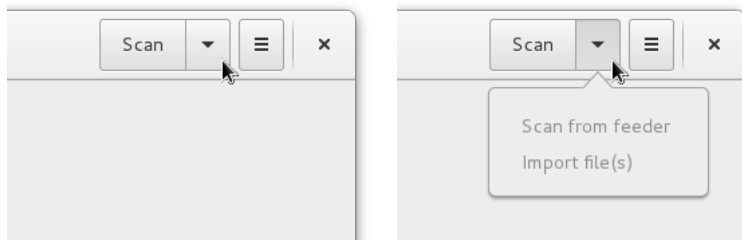


Figure 8: Import option

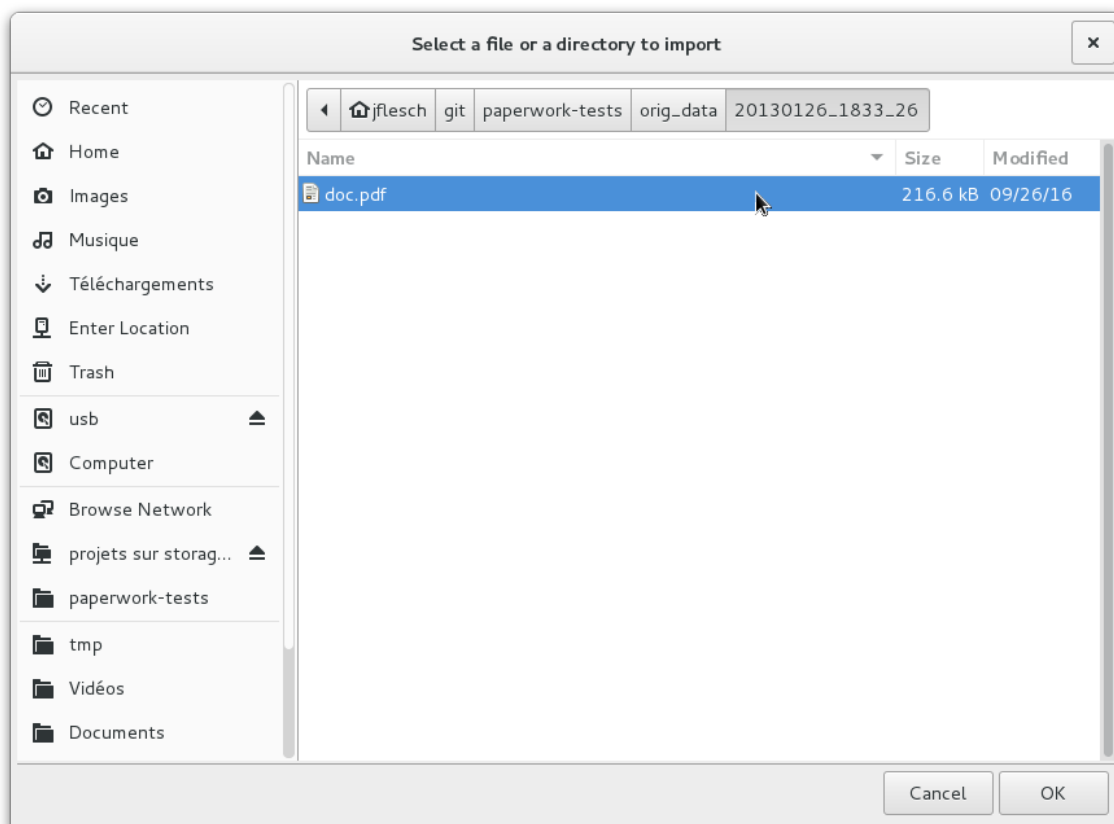


Figure 9: Select a file or a folder

2.1 Images

Paperwork supports a lot of file formats. It supports JPEG, PNG, GIF, BMP, TIFF, etc. Each image is considered as a page. Currently, you can only import one file at a time.

Images are always appended to the document currently opened. Simply select an empty document ("New document") to create a new document while importing.

OCR is always run on imported images. If the imported image is the first page of a new document, Paperwork will automatically apply documents labels.

Note that Paperwork is a document manager. While it can, it is not designed to handle images with only very little text or photos. Automatic labeling will not work correctly on such documents.

The OCR (Tesseract) works very well with black text on white background. Automatic labeling uses recognized text and requires as many keywords on the first page as possible.

2.2 PDF

Each PDF is always considered as a whole document. They are never appended to existing document. They are copied as is in the work directory and are never modified by Paperwork (just moved and renamed).

Paperwork will look for pages with no text attached. On those pages, it will automatically run OCR. Once all the pages have been examined, it will automatically apply document labels. Note that this process may take a few minutes for big PDFs files.

If the PDF is already part of your documents, Paperwork will simply ignore it.

2.3 Many PDFs in one shot

If you import a folder, Paperwork will browse this folder and look for PDFs to import. Already-imported PDFs are simply ignored. Folder is browsed recursively (all the folders inside the folder are also examined).

3 Labels

3.1 Creating new labels

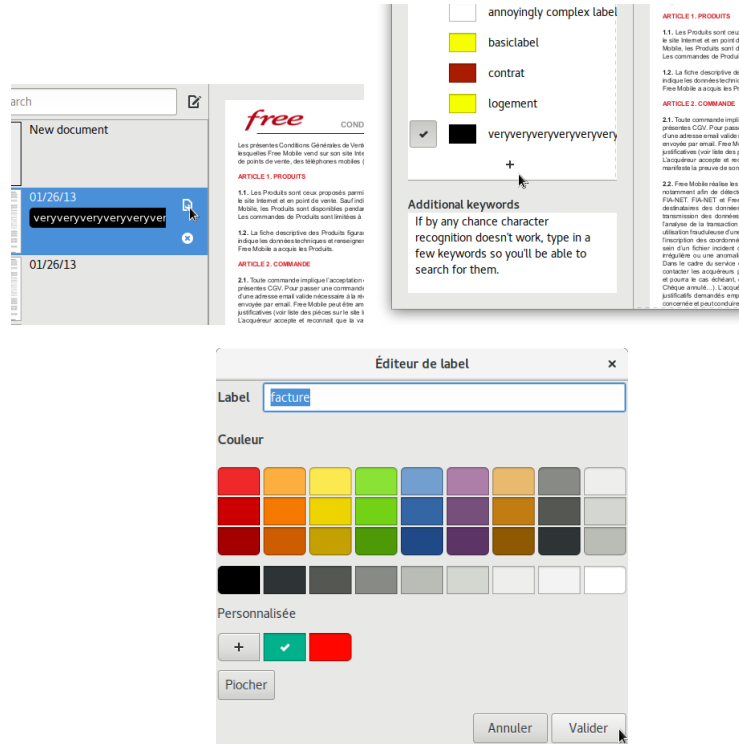


Figure 10: Creating a new label

3.2 Setting labels on documents

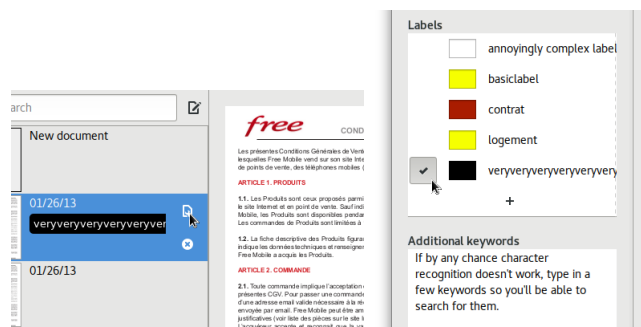


Figure 11: Selecting a label

3.3 Modifying a label

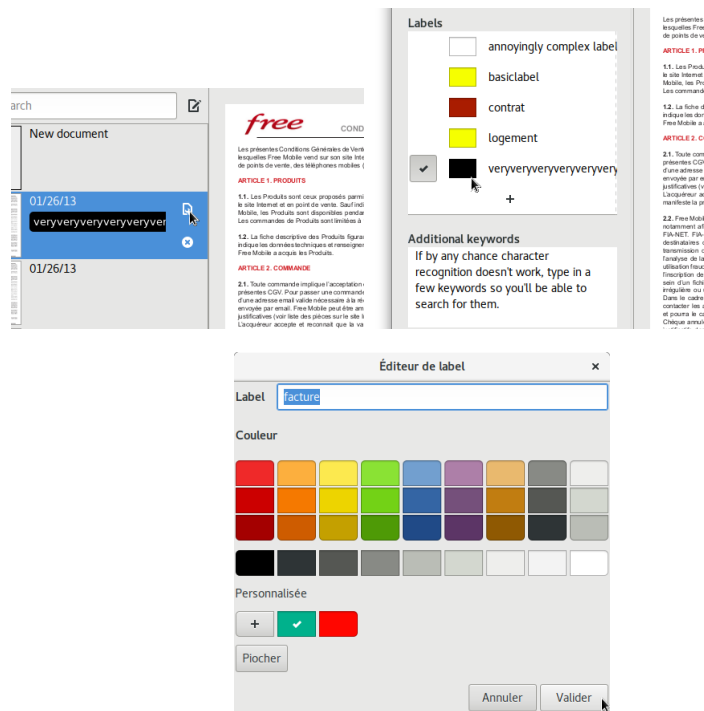


Figure 12: Modifying a label

3.4 Deleting a label



Figure 13: Deleting a label

4 Searching

4.1 Simple search

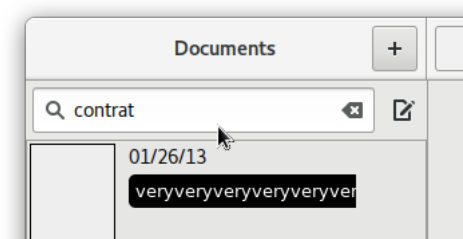


Figure 14: Searching

4.2 Advanced search

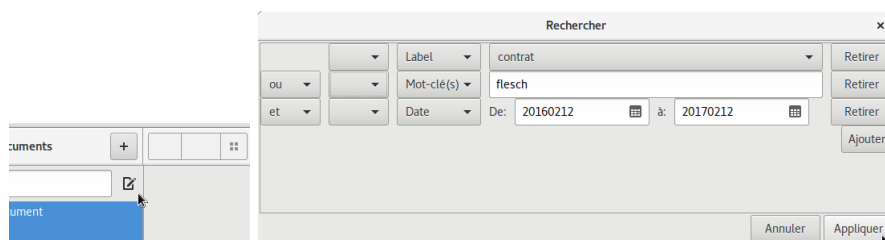


Figure 15: Advanced search

5 Viewing

5.1 View pages as grid

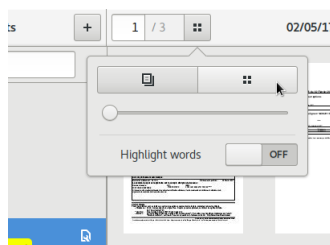


Figure 16: Switch to grid layout

5.2 View pages as list

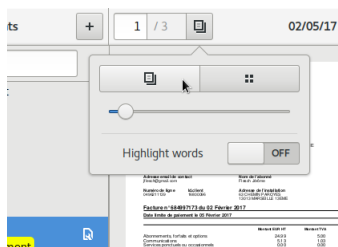


Figure 17: Switch to list layout

5.3 Zoom level

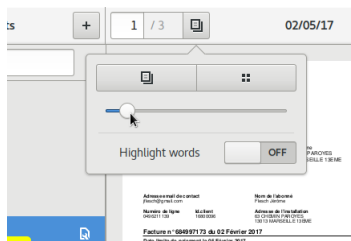


Figure 18: Changing zoom level

6 Exporting

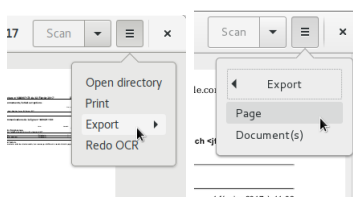


Figure 19: Exporting

6.1 Document

6.2 Page

7 Printing

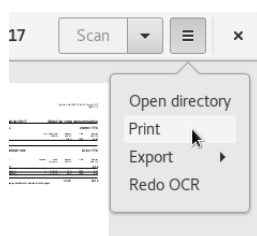


Figure 20: Printing

8 Copying text

9 Editing pages

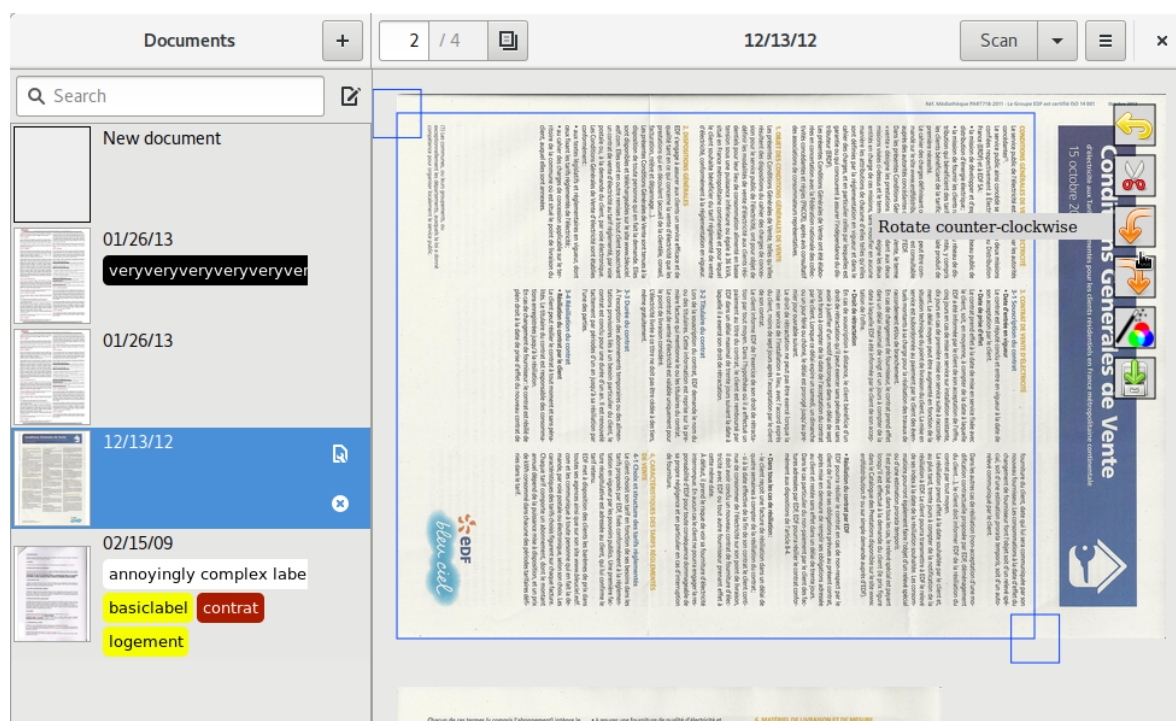


Figure 21: Page editing

10 Moving pages

10.1 inside a document

10.2 from a document to another

11 Switching to another work directory

Before copying or moving the work directory of Paperwork, please close Paperwork.

When Paperwork starts, one of the first things it does is to look for any change in its current work directory. Therefore, if you moved your work directory, when you will restart Paperwork, since it won't find anything, the document list will be empty.

You must then go in the settings and change the work directory location to the new one.

In the following example, we are switching to a work directory contained in a DropBox's folder :

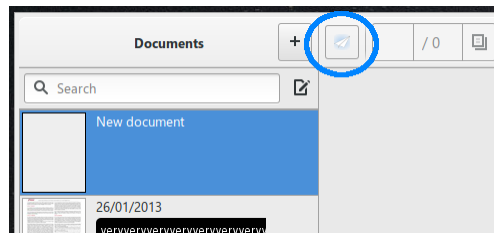


Figure 22: Application menu

Note that, on GNU/Linux, the application menu may be at the top of the screen.

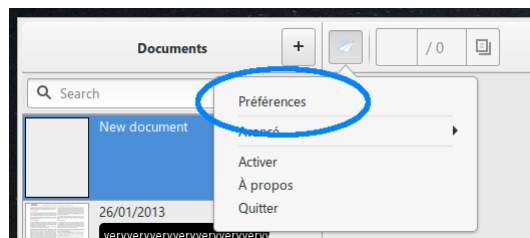


Figure 23: Settings

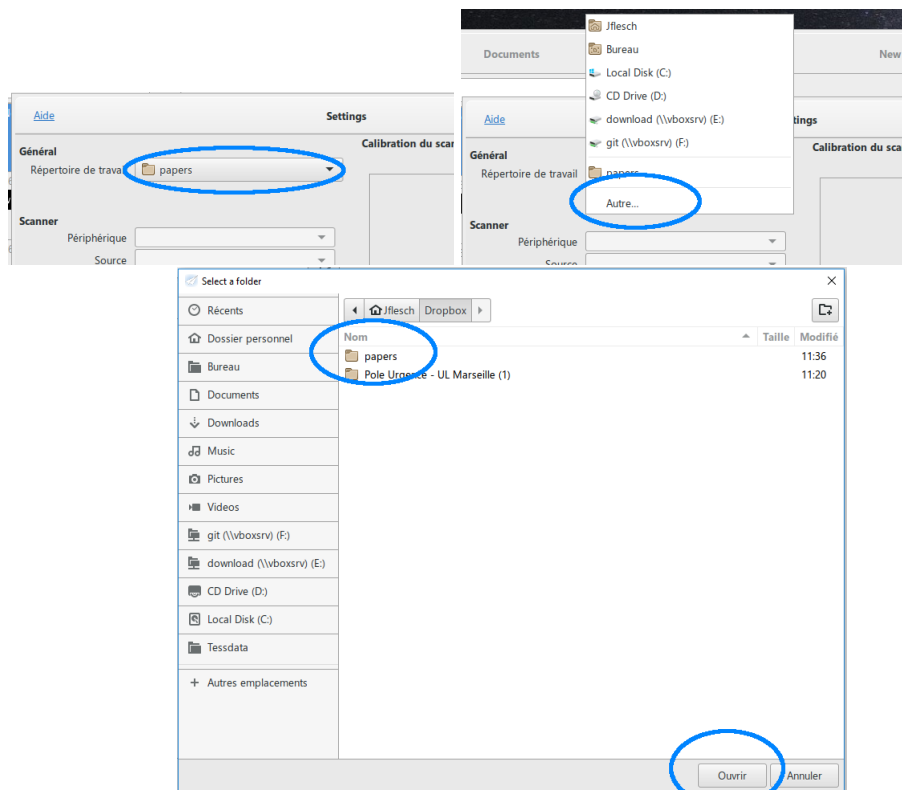


Figure 24: Work directory selection

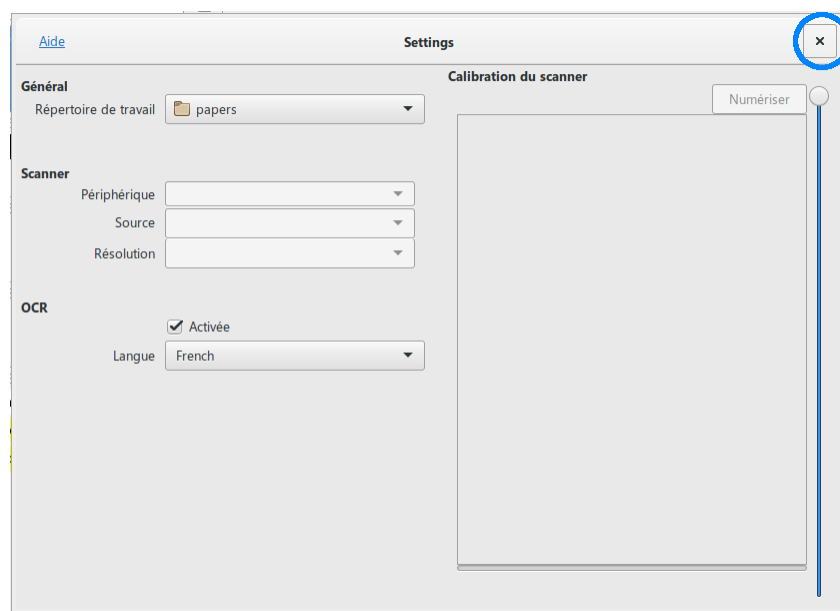


Figure 25: Apply new settings

Paperwork will automatically scan the newly selected work directory, and update its index according to its content.

12 Backup

13 Synchronisation

While Paperwork is a personal document manager, it is not a file synchronization application. However, it is designed to be used with file synchronization applications (Dropbox, OneDrive, Owncloud, SparkleShare, etc).

When you start Paperwork, one of the first things it does is check the content of the work directory. It looks for any changes and updates its document list and index accordingly, automatically.

13.1 USB key / USB drive

This is the simplest way to share documents. Simply copy your work directory to an USB key, tell Paperwork to use it, and you're done.

Beware: You should backup your USB key from time to time on another one.

13.2 File Synchronization applications

Those applications synchronize a local directory with a remote server (or cloud). All the changes you do in your folder are applied on the server. All the changes applied on the servers are applied to the computers that connect to it. The server can belong to you or to someone else (usually a company).

Beware: If you choose to host your documents on someone else server (DropBox, OneDrive, etc), they can access all your documents. Paperwork does not cipher them.

Paperwork is tested daily with SparkleShare. While this is not the easiest one to use, Sparkleshare let you host your files yourself. Using DropBox or OneDrive can make sense if you're sharing not-so-confidential documents with others (associations, etc).

13.2.1 DropBox

Here we are detailing the process to use DropBox, but it is similar for other file synchronization applications.

First, you must copy or move your work directory inside the DropBox folder (please stop Paperwork before):

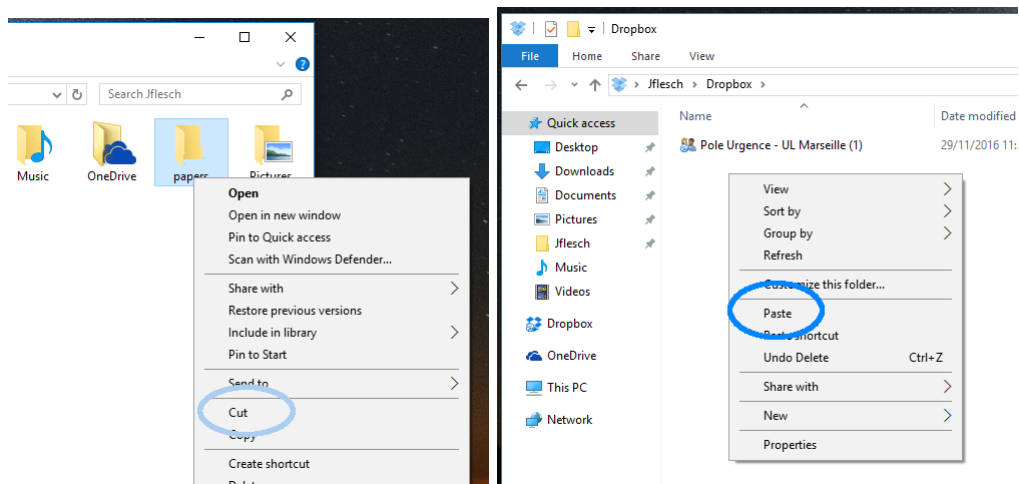


Figure 26: Cut and paste paperwork work directory into the shared directory

Then you must tell Paperwork to use this new work directory.

13.2.2 Shared folder

If all your computers are on the same network, you can share your work directory. However, be really careful regarding permissions. Being too permissive could let a pirate access all your personal documents ! And setting them correctly is tricky.

Beware: While file synchronization applications usually maintain an historic, shared folders do not. You should do backups from time to time.

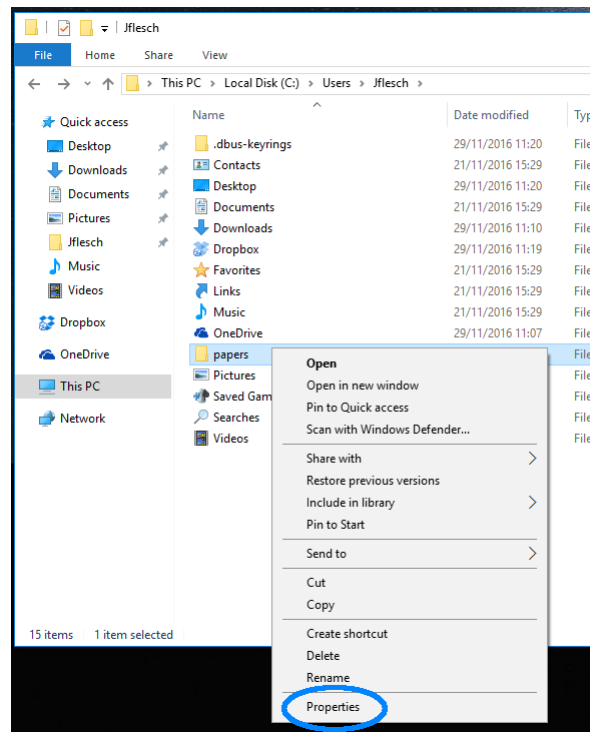


Figure 27: Go to the properties of the work directory

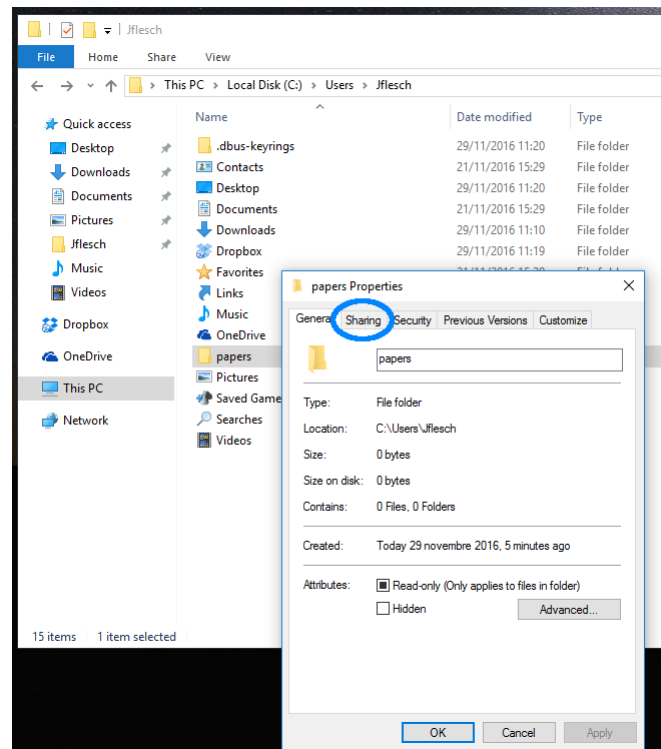


Figure 28: Go to the tab "Sharing"

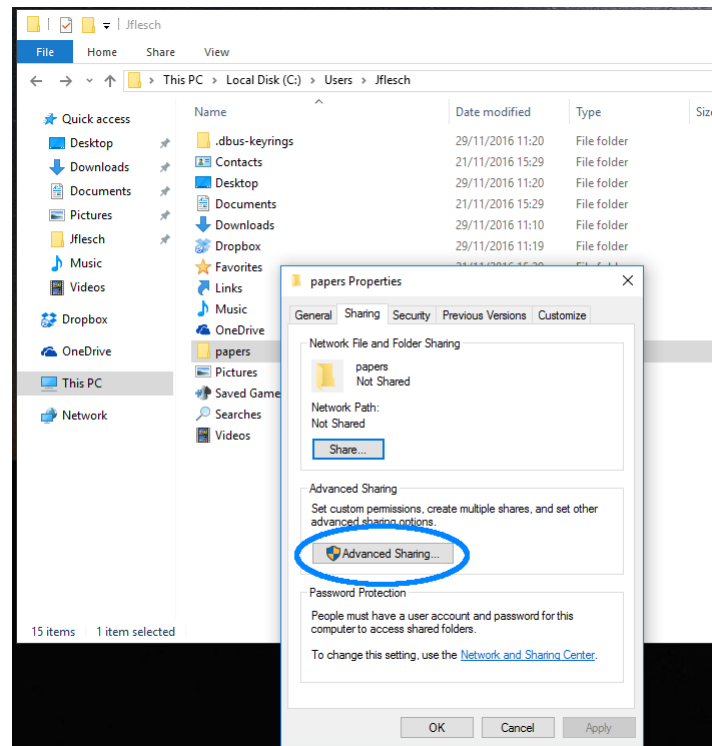


Figure 29: Go to “Advanced Sharing”

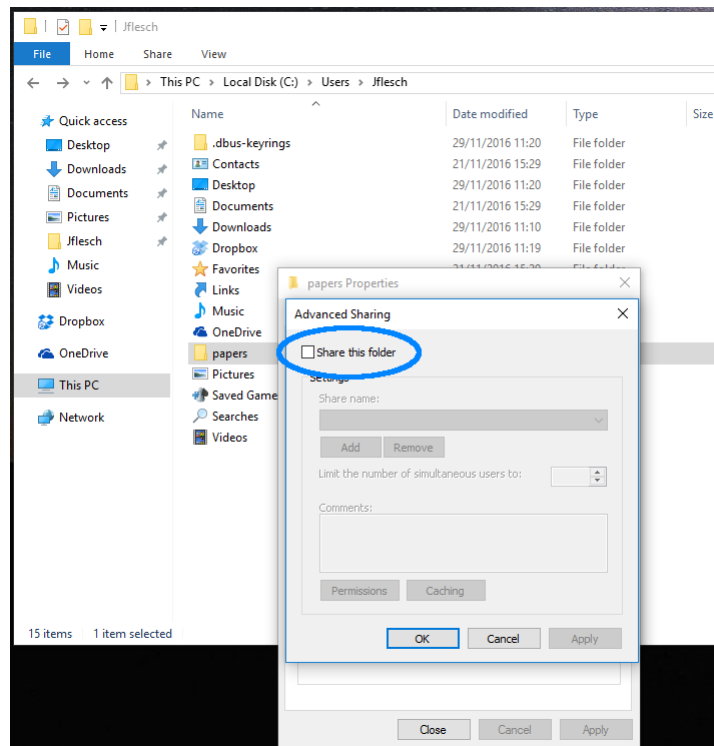


Figure 30: Check “Share this folder”

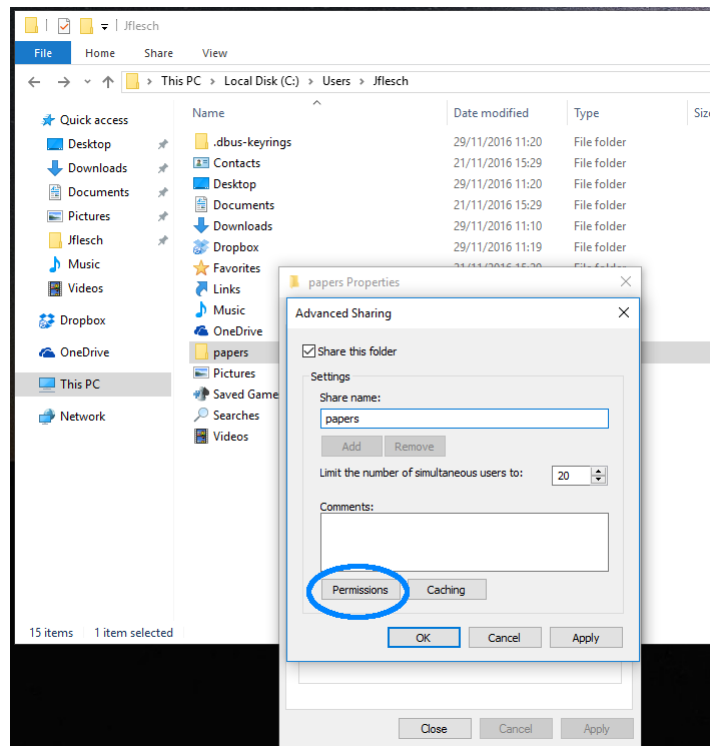


Figure 31: Go to the permissions

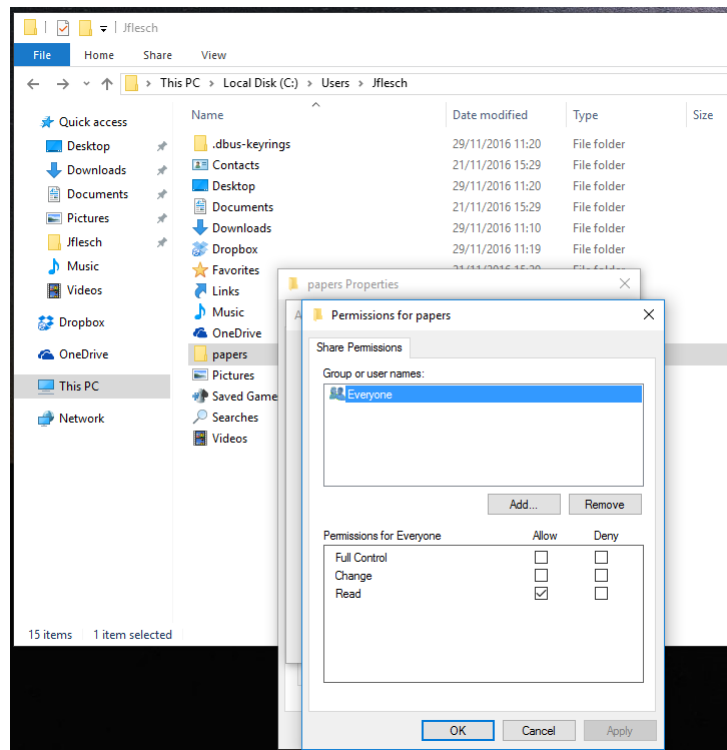


Figure 32: Set the permissions as you wish

Here are the instructions for Microsoft Windows:

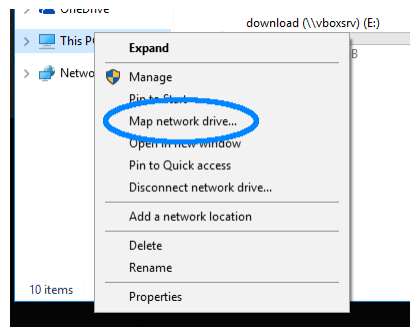


Figure 33: You must map a network drive

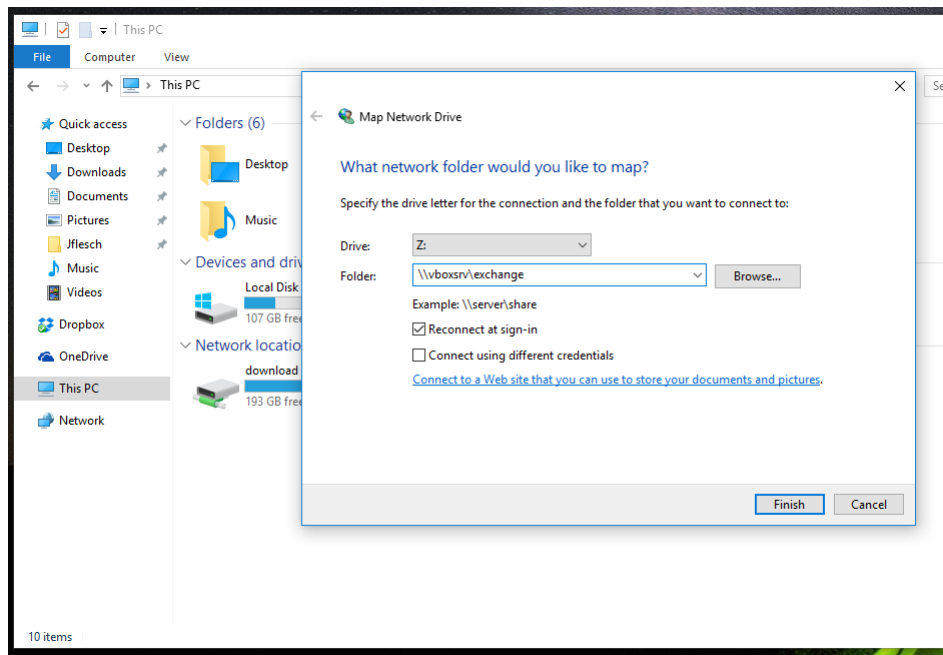


Figure 34: Network drive setup

On the client side, you must map the shared folder to a drive.

14 Encryption

14.1 Windows

14.2 GNU/Linux

GNU/Linux distributions include many tools to encrypt whole directories.

With Paperwork, there are 2 directories that should be encrypted to protect your privacy:

- Your work directory (by default `~/papers`, can be changed in the settings)
- The cache directory (`~/.local/share/paperwork`, cannot be changed) (it contains index files from which the content of your documents could be partially recovered)

14.2.1 Ecryptfs

On GNU/Linux Debian and Ubuntu, you can easily create a directory Private in your home directory. This directory will be encrypted using the password you use to connect when you start your computer. Just type `ecryptfs-setup-private` in a terminal to create it. You have to logout/login again. You can then put the work directory of Paperwork in it.

Once the directory has been created, you can also store Paperwork cache in it:

```
$ mv ~/.local/share/paperwork ~/Private/paperwork_cache
$ ln -s ~/Private/paperwork_cache ~/.local/share/paperwork
```

14.2.2 Encfs

Encfs can also be used to create encrypted directories easily. However, beware that Encfs seems to have some security weaknesses.

```
$ encfs ~/.papers ~/papers
```

15 Advanced use and information

15.1 Redo OCR

15.1.1 On all the documents

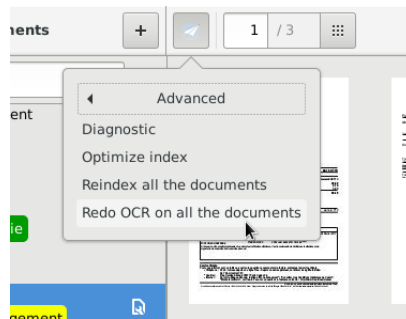


Figure 35: Redo OCR on all the documents

15.1.2 On one document

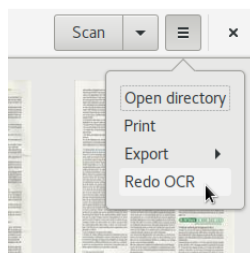


Figure 36: Redo OCR on one document

15.2 Highlight all words

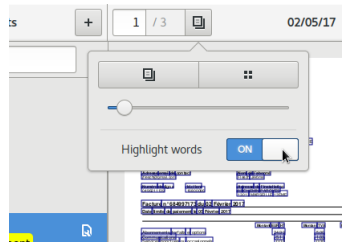


Figure 37: Highlight all words

15.3 Keyboard shortcuts

- Ctrl+E
- Ctrl+N
- PageUp
- PageDown
- Ctrl+PageUp
- Ctrl+PageDown
- Shift+MouseButton on a document

15.4 Paperwork's files locations

By default:

- Configuration : `~/.config/paperwork.conf`
- Index : `~/.local/share/paperwork`
- Documents : `~/papers`

(same paths are used on Windows ; `~` = `C:\Users[login]` ; folders are hidden)

The index is always updated according based on the documents in the work directory. When Paperwork starts, the modification time of each file is used to detect changes on the documents.

15.5 Work directory layout

`workdir|rootdir` = `~/papers` (by default)

15.5.1 Global organisation

In the work directory, you have folders, one per document.

The folder names are (usually) the scan/import date of the document: YYYYMMDD_hhmm_ss[_<idx>]. The suffix 'idx' is optional and is just a number added in case of name collision.

In every folder you have:

- For image documents:
 - paper.<X>.jpg : A page in JPG format (X starts at 1)
 - paper.<X>.words (optional) : A hOCR file, containing all the words found on the page using the OCR (optional, but required for indexing ; can be regenerated with the options "Redo OCR (...)").
 - paper.<X>.thumb.jpg (optional, generated automatically) : A thumbnail version of the page (faster to load) labels (optional) : a text file containing the labels applied on this document
 - extra.txt (optional) : extra keywords added by the user
- For PDF documents:
 - doc.pdf : the document labels (optional) : a text file containing the labels applied on this document
 - extra.txt (optional) : extra keywords added by the user
 - paper.<X>.words (optional) : A hOCR file, containing all the words found on the page using the OCR. Some PDF contains crap instead of the real text, so running the OCR on them can sometimes be useful.

Here is an example a work directory organisation:

```
$ find ~/papers
/home/jflesch/papers
/home/jflesch/papers/20130505_1518_00
/home/jflesch/papers/20130505_1518_00/paper.1.jpg
/home/jflesch/papers/20130505_1518_00/paper.1.thumb.jpg
/home/jflesch/papers/20130505_1518_00/paper.1.words
/home/jflesch/papers/20130505_1518_00/paper.2.jpg
/home/jflesch/papers/20130505_1518_00/paper.2.thumb.jpg
/home/jflesch/papers/20130505_1518_00/paper.2.words
/home/jflesch/papers/20130505_1518_00/paper.3.jpg
/home/jflesch/papers/20130505_1518_00/paper.3.thumb.jpg
/home/jflesch/papers/20130505_1518_00/paper.3.words
/home/jflesch/papers/20130505_1518_00/labels
/home/jflesch/papers/20110726_0000_01f
/home/jflesch/papers/20110726_0000_01/paper.1.jpg
/home/jflesch/papers/20110726_0000_01/paper.1.thumb.jpg
/home/jflesch/papers/20110726_0000_01/paper.1.words
/home/jflesch/papers/20110726_0000_01/paper.2.jpg
/home/jflesch/papers/20110726_0000_01/paper.2.thumb.jpg
/home/jflesch/papers/20110726_0000_01/paper.2.words
/home/jflesch/papers/20110726_0000_01/extra.txt
```

```
/home/jflesch/papers/20130106_1309_44  
/home/jflesch/papers/20130106_1309_44/doc.pdf  
/home/jflesch/papers/20130106_1309_44/paper.1.words  
/home/jflesch/papers/20130106_1309_44/paper.2.words  
/home/jflesch/papers/20130106_1309_44/labels  
/home/jflesch/papers/20130106_1309_44/extra.txt
```

15.5.2 hOCR files

With Tesseract, the hOCR file can be obtained with following command:

```
tesseract paper.<X>.jpg paper.<X> -l <lang> hocr && mv paper.<X>.html paper.<X>.words
```

For example:

```
tesseract paper.1.jpg paper.1 -l fra hocr && mv paper.1.html paper.1.words
```

15.5.3 Label files

Here is an example of content of a label file:

```
facture,#0000b1588c61 logement,#f6b6ffff0000
```

It's always [label],[color]. For a same label, the color should always be the same.

15.6 Statistics

You can get various statistics regarding your documents. Just have a look at the diagnostic output. Statistics are close to the end of the output.

16 Getting support / reporting issues

16.1 Diagnostic dialog

16.2 Github issue tracker

16.3 Mailing-list

17 Uninstalling

Paperwork can be uninstalled. Uninstalling Paperwork *won't* remove your work directory or documents.

17.1 Windows

17.2 GNU/Linux

If you installed Paperwork manually:

```
sudo pip uninstall paperwork  
sudo pip uninstall pyocr  
sudo pip uninstall pyinsane
```

(it's python-pip on some systems)

If you installed many versions of these packages, you may have to run these commands many times.

Note that there are other dependencies installed with Paperwork. However, python-pip can't detect and remove automatically unused dependencies. This is why you should use your distribution package(s) if possible.