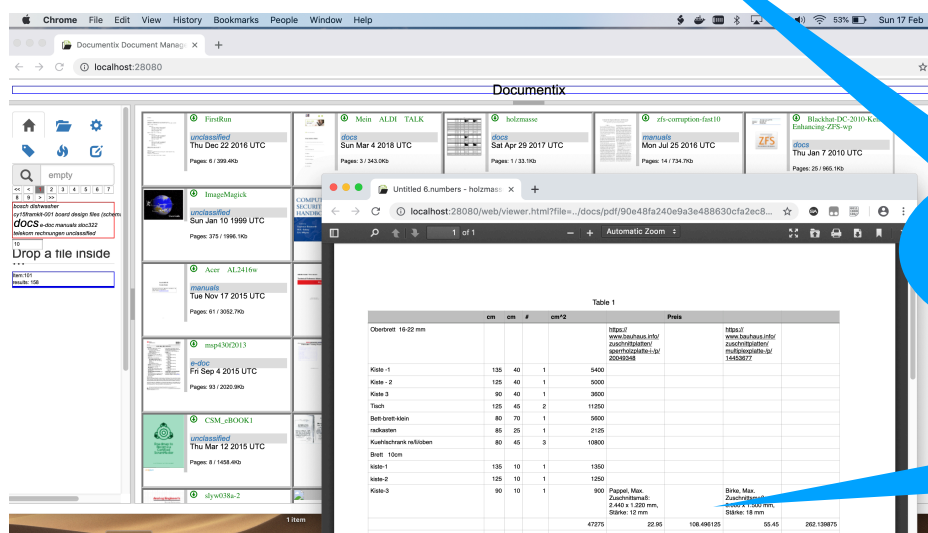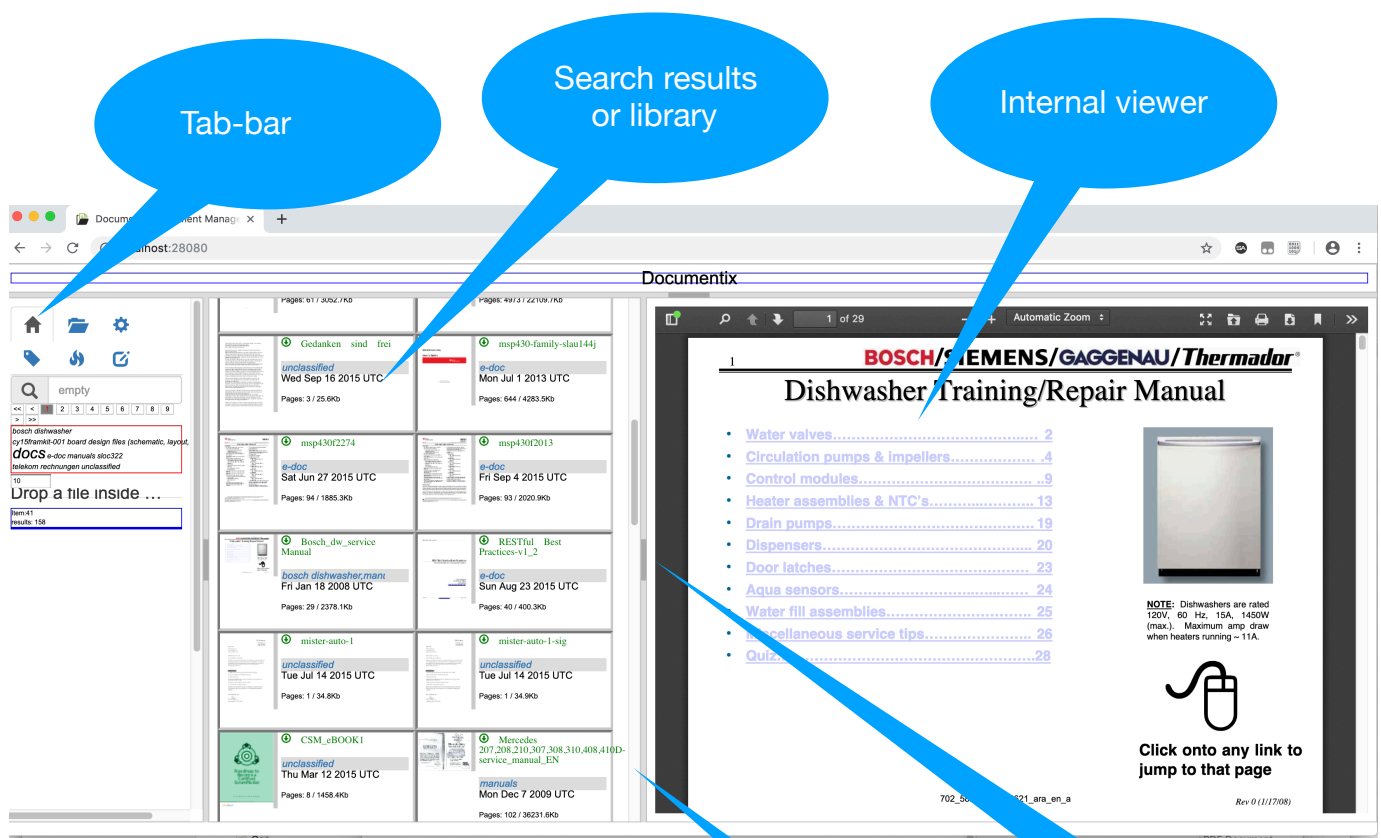# Documentix

## the Document manager

Overview

Documentix quickly searches your collection using a full-text search engine
Documentix lets you classify documents
Documentix automatically classify new documents based on previous training
Documentix automatically OCR's new documents to enable searching
Documentix manages multiple types of documents (pdf/word/ppt/jpeg..)
Documentix presents all archived documents through a web-interface
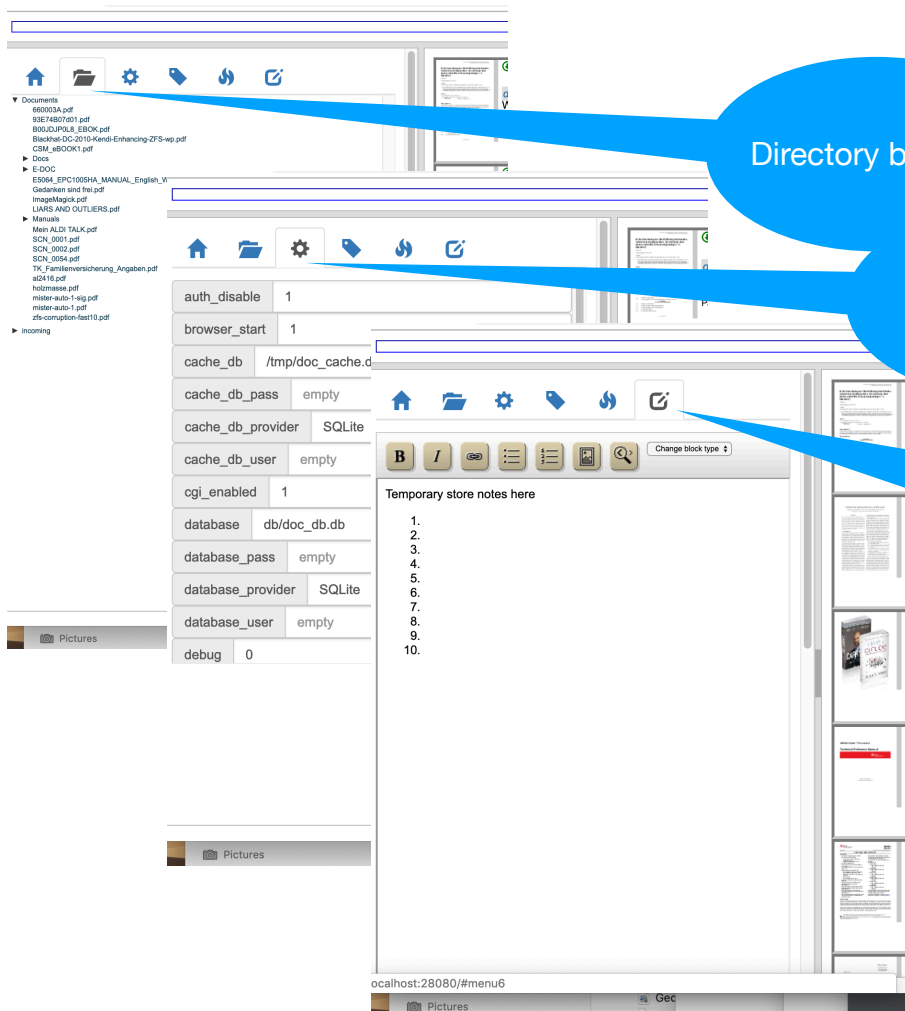Documentix allows easy sharing of documents

## User interface



Tab-bar

Search results or library

Internal viewer

Adjustable layout

Continous scroll bar

External viewer

# Search function



**Download button**

**Document tags**

Lagebericht2014

docs
Thu Dec 11 2014 UTC

Pages: 44 / 1617.0Kb

linux

<< < 1 > >>

*docs* e-doc manuals unclassified

10

Drop a file inside

Item:20

**Search string**

...32

3.3.4

Operation Windigo – **Linux**-Schadpr

32

**Preview of found text**

# Upload new documents

Drag & Drop area

Click to get file-upload dialog

Directory browser

Settings

Scratch-pad

# Drag & Drop Bookmark-icon to capture document & location in scratch-pad

Documentix

bimschv-gesamt
docs
Wed Feb 26 2014 UTC

FirstRun
unclassified
Thu Dec 22 2016 UTC

Mein ALDI TALK
docs
Sun Mar 4 2018 UTC

holzmasse
docs
Sat Apr 29 2017 UTC

Change block type

Temporary store notes here

1. Check this for more info:
2. bimschv-gesamt.pdf(4)
3.
4.
5.
6.
7.
8.
9.

bimschv-gesamt.pdf

localhost:28080/web/viewer.html?file=../docs/pdf/c1b1cdd322ee72101b9950519164d69...

5 of 23        Automatic Zoom

- Seite 4 von 23 -

Ein Service des Bundesministeriums der Justiz und für Verbraucherschutz
in Zusammenarbeit mit der juris GmbH - www.juris.de

9.    Heizöl leicht (Heizöl EL) nach DIN 51603-1, Ausgabe August 2008, und andere leichte Heizöle mit gleichwertiger Qualität sowie Methanol, Ethanol, naturbelassene Pflanzenöle oder Pflanzenölmethylester,

10.   Gase der öffentlichen Gasversorgung, naturbelassenes Erdgas oder Erdölgas mit vergleichbaren Schwefelgehalten sowie Flüssiggas oder Wasserstoff,

11.   Klärgas mit einem Volumengehalt an Schwefelverbindungen bis zu 1 Promille, angegeben als Schwefel, oder Biogas aus der Landwirtschaft,

12.   Koksofengas, Grubengas, Stahlgas, Hochofengas, Raffineriegas und Synthesegas mit einem Volumengehalt an Schwefelverbindungen bis zu 1 Promille, angegeben als Schwefel, sowie

# Internals

Documentix uses SQLite and the full-text-search engine as a database backend.
All documents are stored using the MD5 hash as an index. (Don't tell me it's broken - it does not matter here)

Documents send to the processing engine are converted to PDF'files (if not already). If the PDF does not contain more than a certain number of words, tesseract as an OCR engine is started on the document and a searchable PDF is created.

An interface to a Canon-mx870 exists but is not supported in the docker image.

It should be easy to create a curl command-line to upload new scans like:
> curl  -F file=@o-page-14.pdf -H x-file-name:MyPage.pdf http://localhost:28080/upload

Which will upload the file "o-page-14.pdf" and name it on the server "MyPage.pdf"

Retrieving files is also possible:

If you know the MD5 has of a file.

Lets assume the MD5="009914b6b6a4cb4a37f2fc5bdbc561ca" is.
And the server is reachable under: http://localhost:28080

Then:

http://localhost:28080/docs/raw/009914b6b6a4cb4a37f2fc5bdbc561ca/MyPage
    Will return the original document as it was uploaded


http://localhost:28080/docs/pdf/009914b6b6a4cb4a37f2fc5bdbc561ca/MyPage
    Will get the file as a PDF (converted if necessary including OCR)


http://localhost:28080/web/viewer.html?file=../docs/pdf/009914b6b6a4cb4a37f2fc5bdbc561ca/MyPage
    Will show the file as a PDF in the html viewer.

http://localhost:28080/docs/ico/009914b6b6a4cb4a37f2fc5bdbc561ca/MyPage
    Will return the original document as it was uploaded