



Strigi in KDE4
the power of indices

Jos van den Oever



History of free desktop search





1996: KFind

2001: KFileMetaInfo

2005: start of Kat

aKademy 2005: Kat and Tenor hype

aKademy 2006: Nepomuk and Strigi are presented

Now

Nepomuk

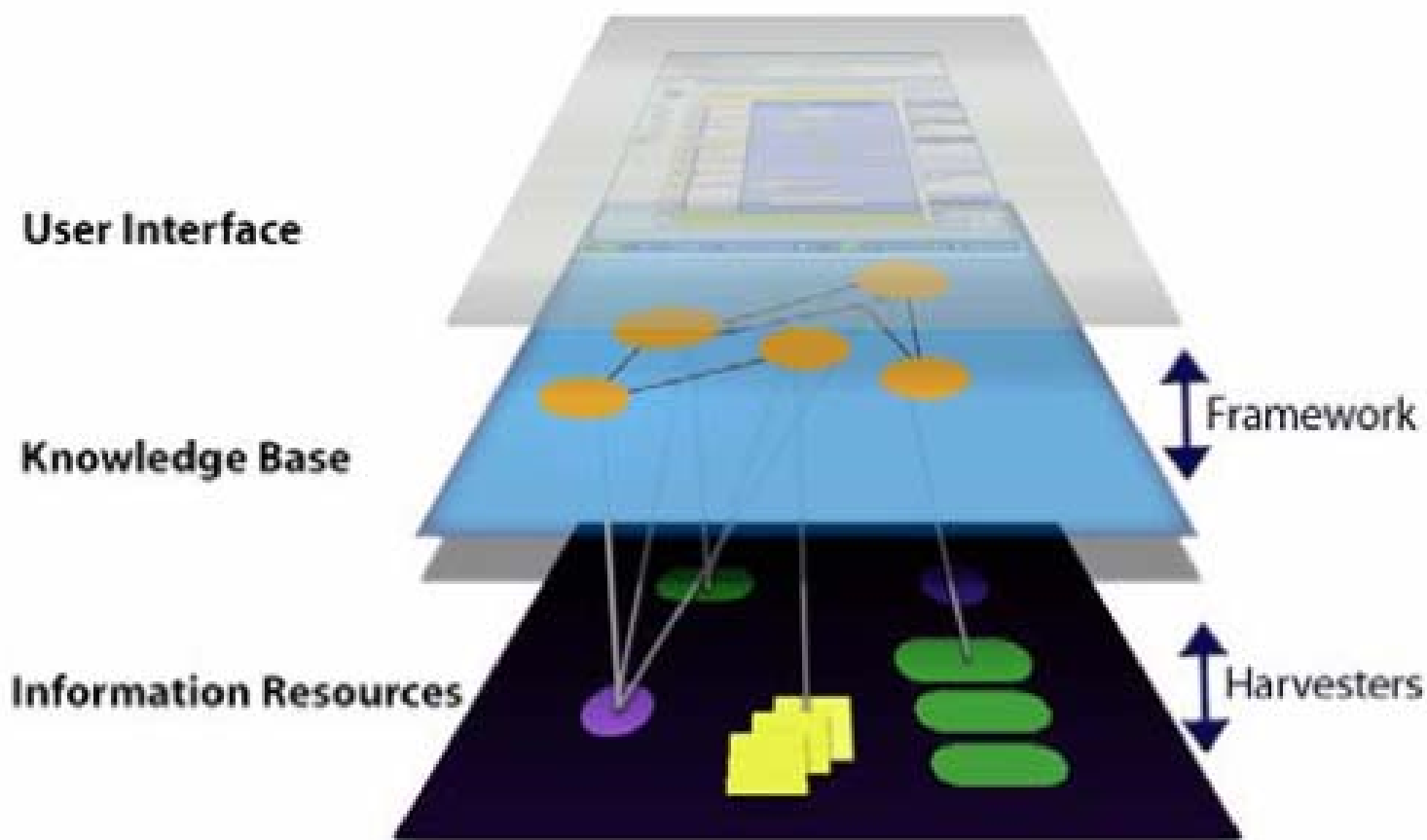
semantic storage
and standards

Strigi

data extraction,
indexing, search

Xesam

freedesktop.org
search standard





libstreams

- efficient streaming access to file contents
- universal API to different formats

libstreamanalyzer

- analysis of libstreams streams with many parallel analyzers
- storage and retrieval over abstract interface



Reading nested files

*.gz	zcat
*.bz2	bzcat
*.tar	tar
*.zip, *.tar, *.jwe, openoffice files	unzip
email	mail client
email attachment	mail client
*.pdf (?)	?
*.deb, *.ar, static libs	ar
*.cpio	cpio
*.rpm	rpm2cpio + cpio

many formats, many tools, many interfaces



Common API for nested files

Can we use kio or vfs?



tar:/home/me/data.tar/file1



zip:/another.zip/example.txt

zip:/

tar:/

gz:/

rpm:/

deb:/

disadvantages:

- user has to figure out what kio or vfs is required

solution:

- make a clever kio/vfs that understands all

commonapi:/

alternative: fuse



```
tar:/home/me/data.tar/file1.zip#zip:example.txt
```

“None of the chained uri stuff (tar/zip/etc)

really work, and never did.”

Alexander Larsson,

Oct 2005 to gnome-vfs-list@gnome.org

“Bug 73821: Please "unchain" kioslaves.

Browsing a zip inside a zip should work.”

KDE bug since Jan 2004

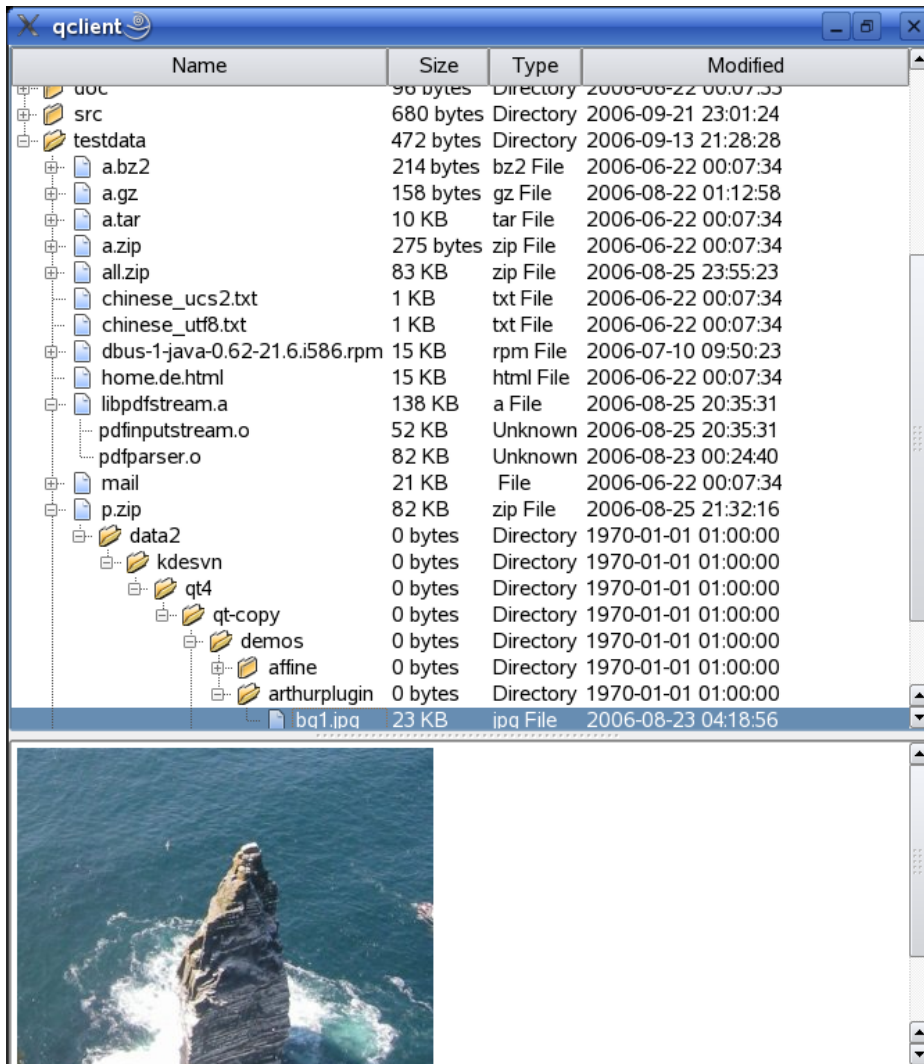


StreamBase and SubStreamProvider

```
class StreamBase {
    virtual int32_t read(const char** data, int32_t min, int32_t max) = 0;
    int64_t reset(int64_t newpos) = 0;
};
```

```
void
readdemo() {
    int32_t nread;
    const char* data;
    nread = stream->read(data, 1, 0); // read at least 1 byte
    stream->reset(0); // reset to start of stream
    nread = stream->read(data, 3, 3); // read exactly 3 bytes
}
```

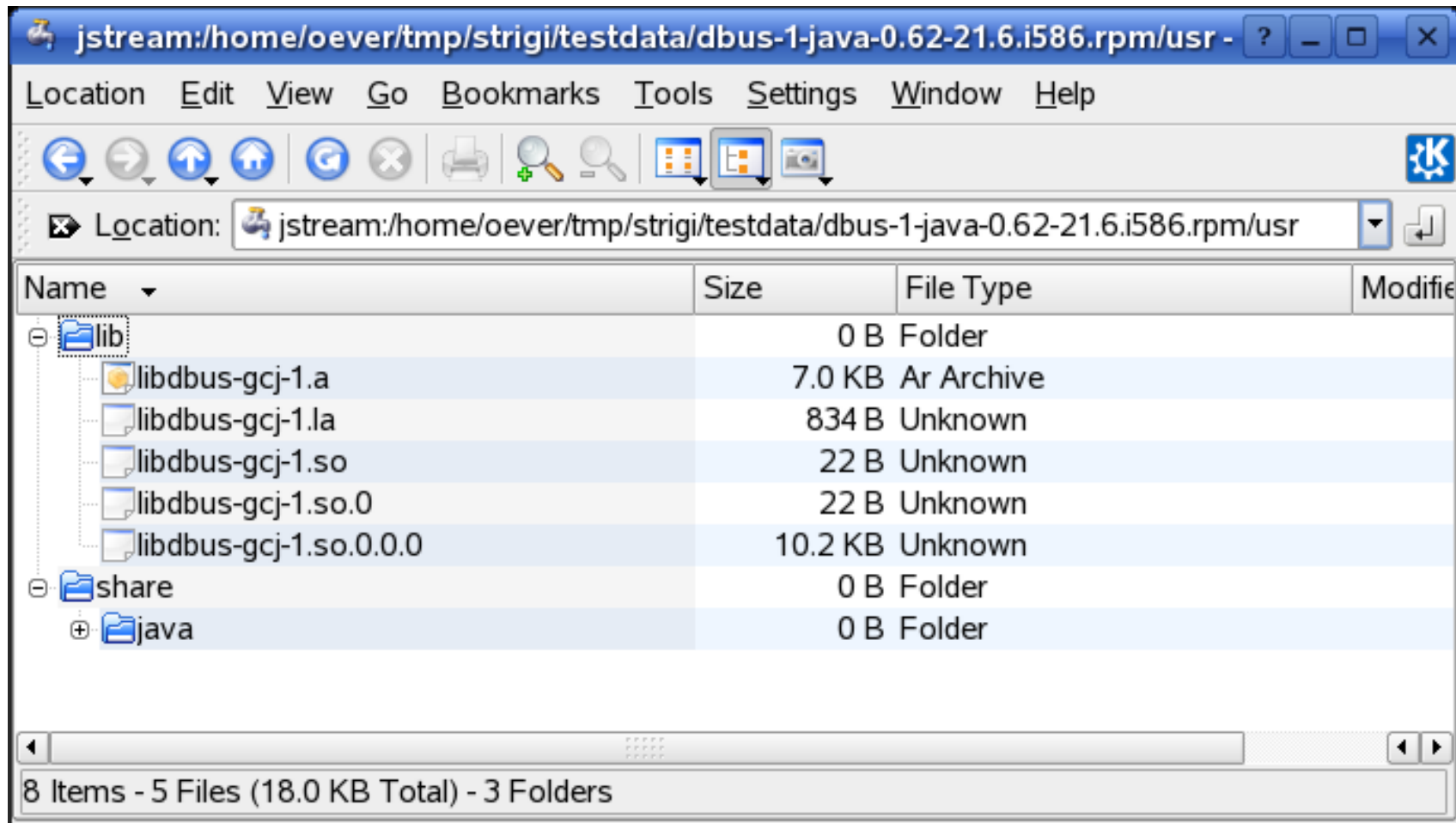
```
class SubStreamProvider {
    virtual int32_t read(const char** data, int32_t min, int32_t max) = 0;
    virtual int64_t reset(int64_t newpos) = 0;
};
```



add read access to archive
formats by adding only one
line of code:

ArchiveEngineHandler engine;

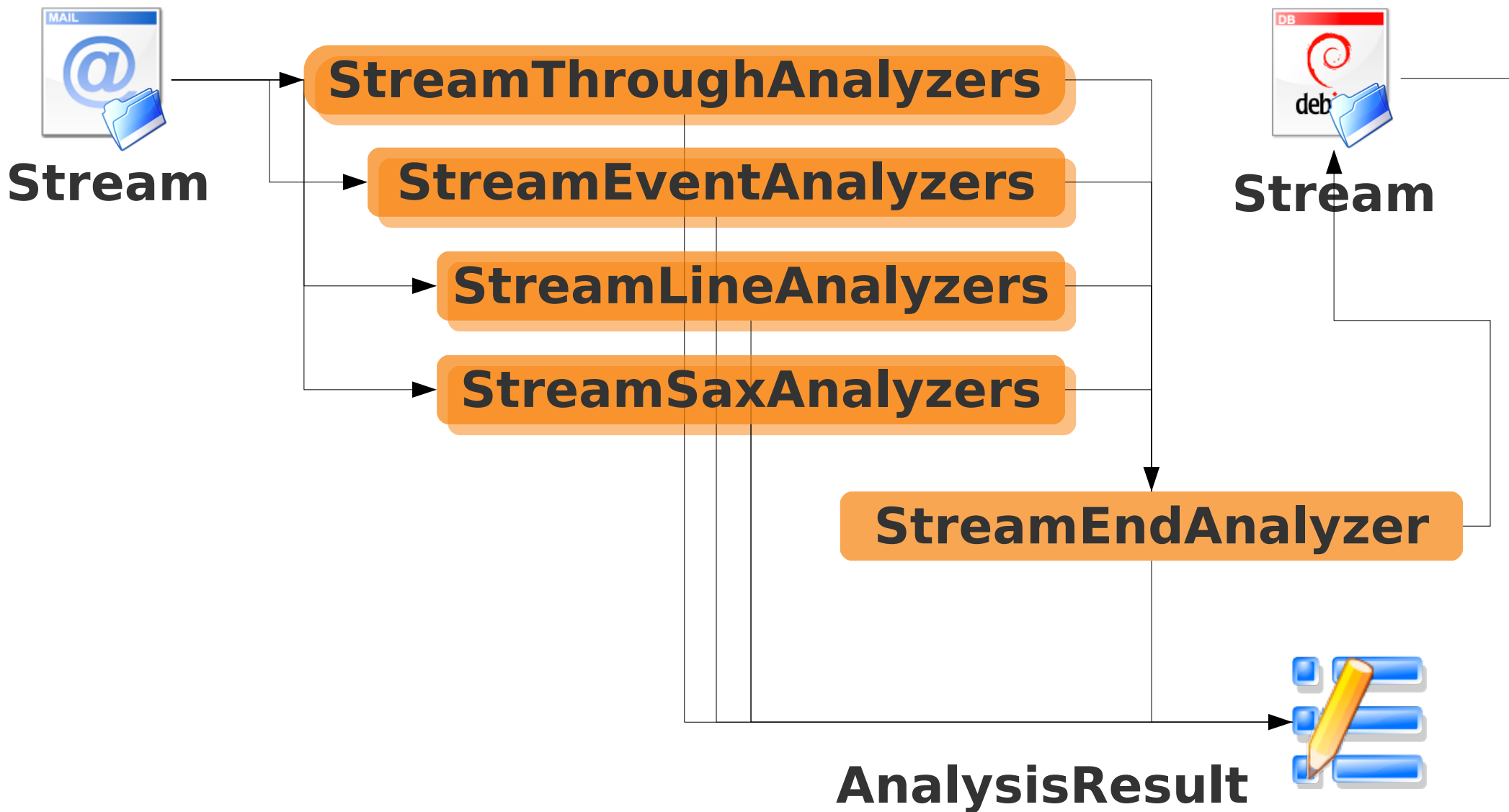
Class that comes with Strigi that uses
QabstractFileEngine to give Qt
applications transparent access to a
custom filesystem.





directory | file





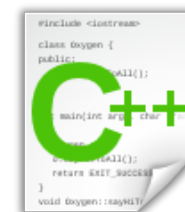


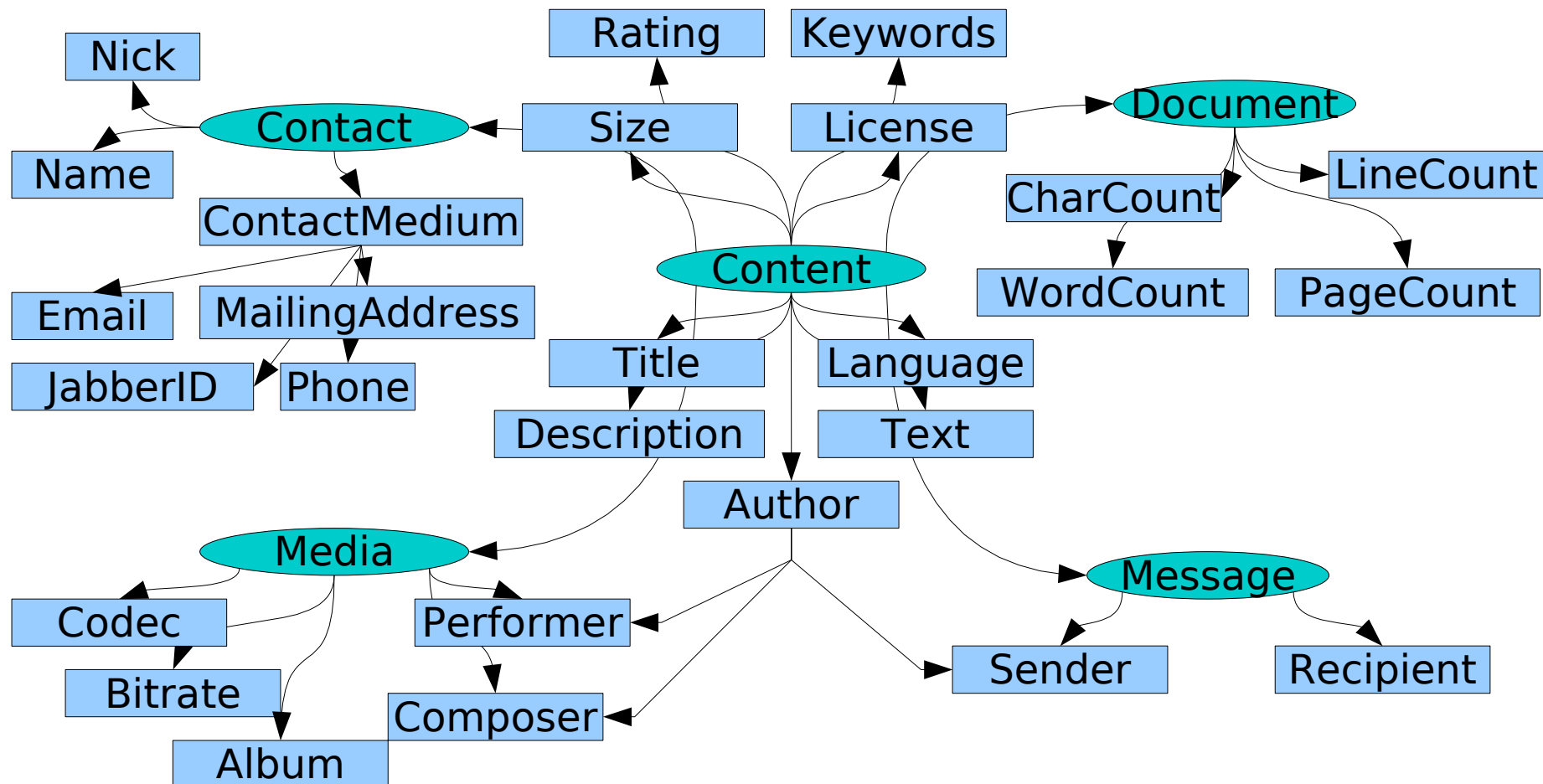
```
class RegExLineAnalyzerFactory : public LineAnalyzerFactory {  
    StreamLineAnalyzer* newInstance() const;  
};
```

```
class RegExLineAnalyzer : public StreamLineAnalyzer {  
public:  
    void startAnalysis(Strigi::AnalysisResult*);  
    void handleLine(const char* data, uint32_t length);  
    void endAnalysis();  
    bool isReadyWithStream();  
};
```



Selection of file formats







Indexes and Index Management

IndexManager

IndexReader

IndexWriter

Indexes

Glucene

Soprano

SQLite

HyperEstraiier

Xapian

semi-Indexes

KFileMetaInfo

CombinedIndexReader

GreplIndex

`xmlindexer`

`deepfind`

`deepgrep`



strigicmd

create, query, inspect
indexes from the
command line

libstreams **libxml** **libbz2**

libdbus-1 **libclucene**

libz **libstreamanalyzer**

3 MB resident memory

strigidaemon

connection protocols

dbus

unix socket

web service

interfaces

Xesam Live Query

Strigi

implementation

multithreaded queue

configuration

indices



Indexing 10 000 text files (168 MB)

Beagle	2h18	12m
Jindex	3h02	9m
Tracker	3h03	142m
Strigi	0h04	>4m

Source: Comparison of indexers
November, 2006
Michal Pryc, Xusheng Hui
Sun Microsystems



API changed to fit to common ontology

mostly implementation changes

- KFilePlugin changed
 - Strigi<X>Analyzer for reading
 - KFileWritePlugin for writing
- libstreamanalyzer calls many analyzers on each file
- fieldnames changed: ontology is used



<http://nepomuk.semanticdesktop.org>

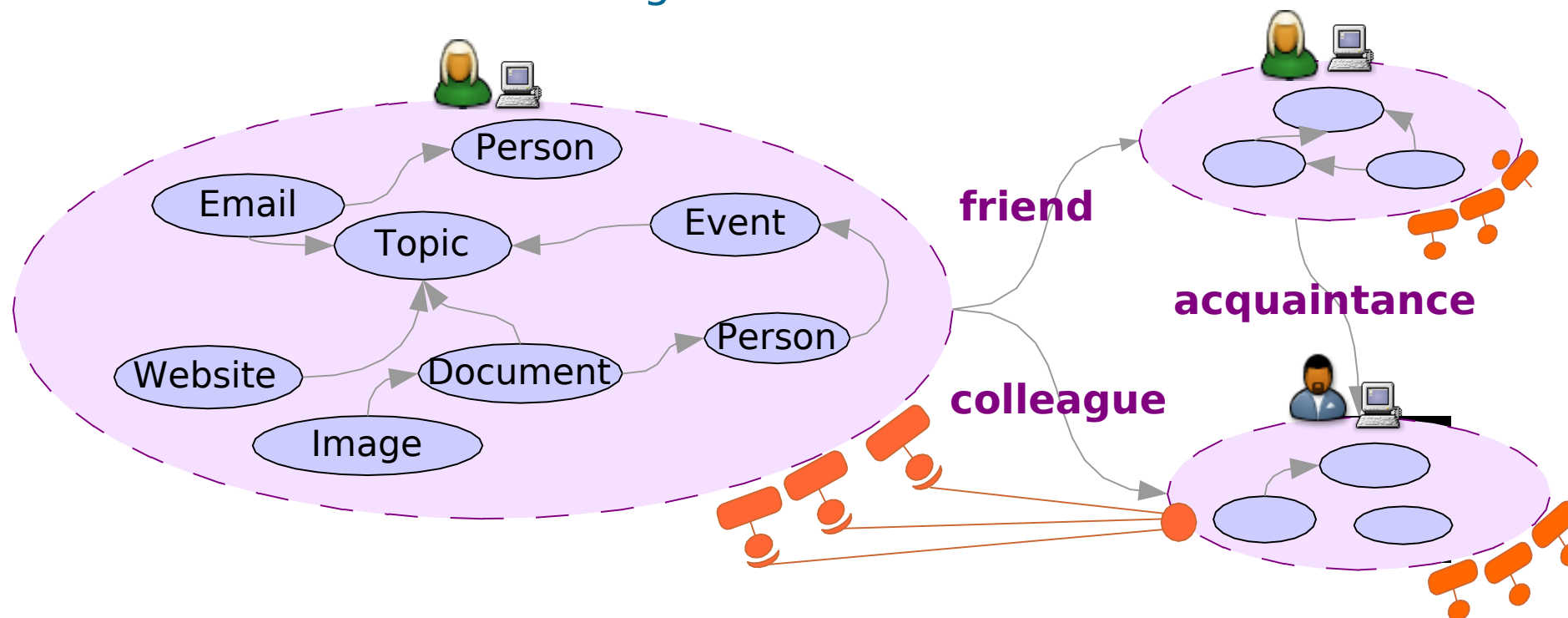


The desktop is a privileged adoption channel for the Semantic Web

Desktop: Help individuals in managing information on the Web/their PC

Semantic: Make content available to automated processing

Social: Enable exchange across individual boundaries



Personal Semantic Web *semantically enlarged intimate supplement to memory*

Social protocols **Social semantic peers and distributed search**



eXtensible Search And Metadata specification

- Dbus API for searching
- fieldnames for standardization

<http://freedesktop.org/wiki/XesamAbout>



Pinot



Nepomuk



Recoll



Strigi



Beagle



Tracker

+ Mikkel Kamstrup Erlandsen



DBus interfaces

- GetHits (in s search, in i num, out aav hits)
- GetHitData (in s search, in ai hit_ids, in as properties, out aav hit_data)

User Query Language

- type:music hendrix

XML Query Language

- `<query><contains><field name="dc:title">
<string>Gödel</string></contains></query>`

Core Ontology



18 chemical formats:

(xyz, vmd, shelx, pdb, mol2, mdl, gaussian, cif, alchemy, cml, ...)

3 streamanalyzers:

(lineanalyzer, saxanalyzer, eventanalyzer)

19 fieldproperties:

(chemistry.inchi,
chemistry.molecular_weight,
chemistry.molecular_formula, ...)

libOpenBabel to
generate InChI

strigi:/?q=chemistry.atom_count:4

search status preferences help about

chemistry.atom_count:4 Found 3 results.

```
<?xml version="1.0"?>
<molecule xmlns="htt
<formula concise="
<identifier version
<basic>1/3N/1/1/3<
</identifier>
<name convention="I
<atomArray>
<atom id="a1" elem
<atom id="a2" elem
<atom id="a3" elem
<atom id="a4" elem
</atomArray>
<bondArray>
<bond atomRefs2="1
```

azane.cml

<?xml version="1.0"?> <molecule
xmlns="http://www.xml-cml.org/schema/cml2/core"
id="CS_azane"> <formula concise=" H 3 N 1 "/>
<identifier version="InChI/1"> <basic>1/
/usr/local/share/chemical-structures/amines/
azane.cml - 1k - XML Document

```
<?xml version="1.0"?>
<molecule xmlns="htt
<formula concise="
<identifier version
<basic>1/2C/1/1/2<
</identifier>
<name convention="I
<atomArray>
<atom id="a1" elem
<atom id="a2" elem
<atom id="a3" elem
<atom id="a4" elem
</atomArray>
<bondArray>
<bond atomRefs2="1
```

acetylene.cml

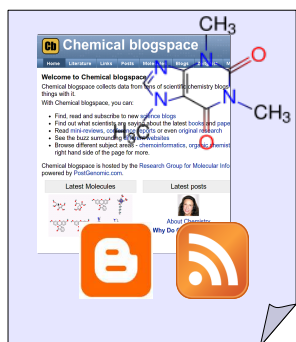
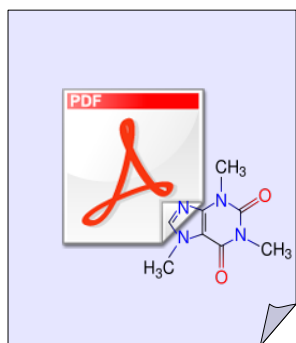
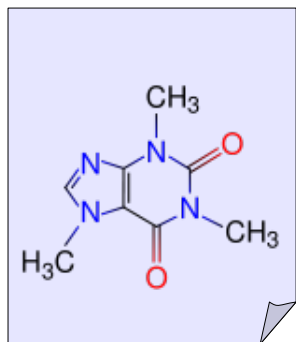
<?xml version="1.0"?> <molecule
xmlns="http://www.xml-cml.org/schema/cml2/core"
id="CS_acetylene"> <formula concise=" C 2 H 2
"/> <identifier version="InChI/1"> <basic>g
/usr/local/share/chemical-structures/alkynes/
acetylene.cml - 1k - XML Document

```
<?xml version="1.0"?>
<molecule xmlns="htt
<formula concise="
<identifier version
<basic>1/2O/1/1/2<
</identifier>
<name convention="I
<atomArray>
<atom id="a1" elem
<atom id="a2" elem
<atom id="a3" elem
<atom id="a4" elem
</atomArray>
<bondArray>
<bond atomRefs2="1
```

formaldehyde.cml

<?xml version="1.0"?> <molecule
xmlns="http://www.xml-cml.org/schema/cml2/core"
id="CS_formaldehyde"> <formula concise=" C 1 H 2
O 1 "/> <identifier version="InChI/1"> <
/usr/local/share/chemical-structures/aldehydes/
formaldehyde.cml - 1k - XML Document

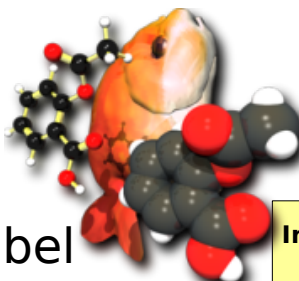
Alexandr Goncarencu, Egon
Willighagen



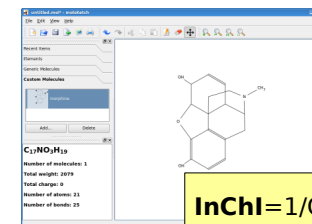
Chemical MIME



Strigi

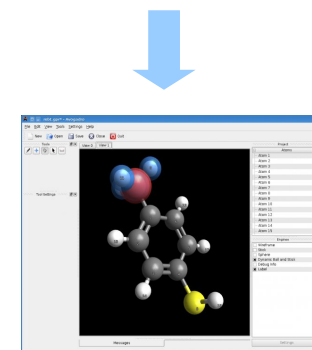


libOpenBabel

CN1C=NC2=C1C(=O)N(C)C(=O)N2C
 InChI=1/C8H10N4O2/
 c1-10-4-9-6-5(10)

CN1C=NC2=C1C(=O)N(C)C(=O)N2C
 InChI=1/C8H10N4O2/
 c1-10-4-9-6-5(10)

molsketch

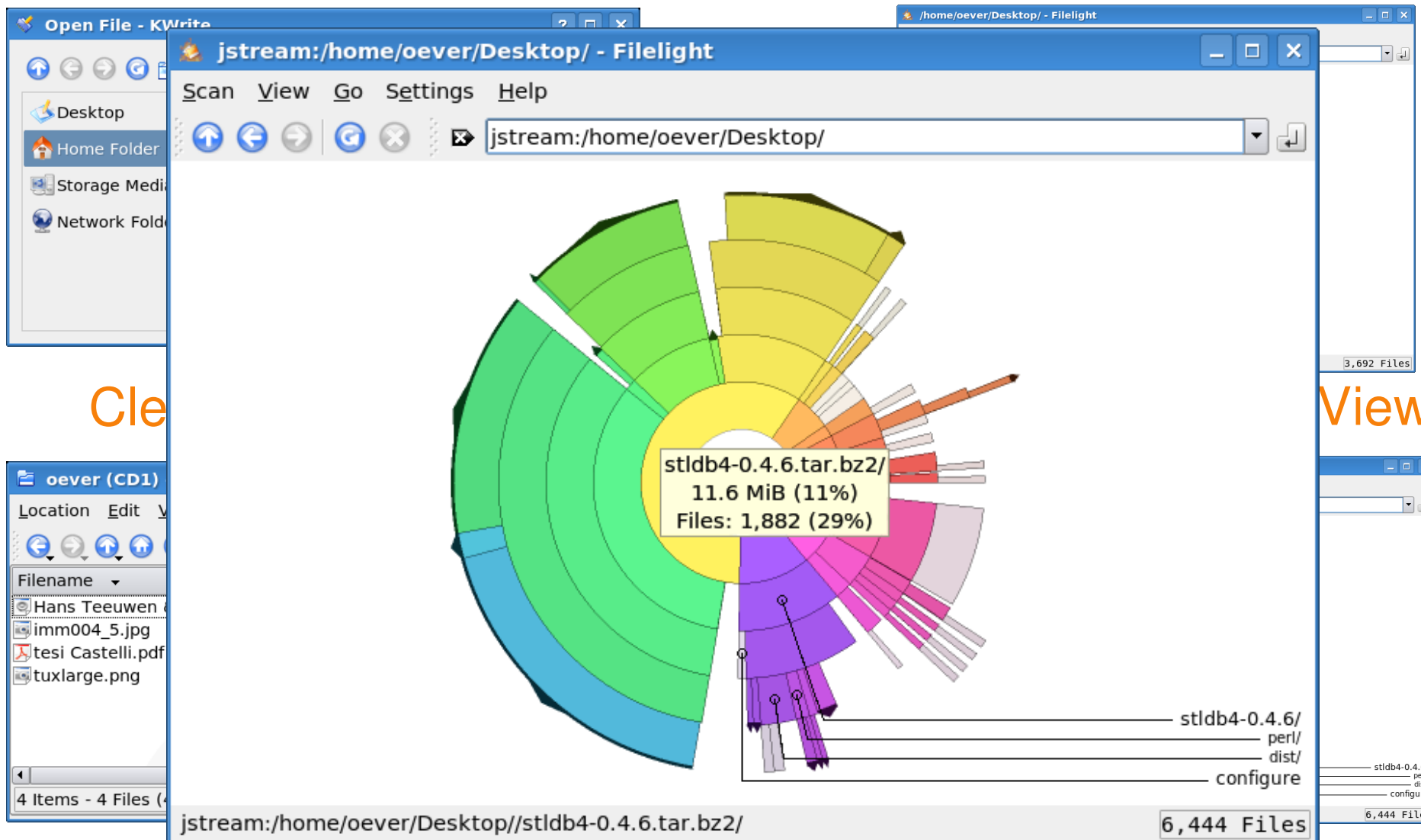
List of search results



Kalzium/Avogadro



File Manager improvements



Cle

View

Clever File Dialog

Universal Radial View



fast stream libraries for reading and analyzing streams

use of modern technologies with a wide consensus

power of a indices to make your applications fast and clever

KDE 4

Nepomuk

semantic storage
and standards

Strigi

data extraction,
indexing, search

Xesam

freedesktop.org
search standard



- + is widely deployed and tested on other platforms
- + has a stable well documented API
- + has a documented API for querying the search daemon
- is closed source software
- uses a proprietary index format
- uses COM for communication
- has a large brand recognition and there will a demand for it
- calls analyzer plugins based on file extension
- has a limited, unexpandable list of categories for files
- identifies files by mtime + uri
- uses wchar_t internally
- is file based
- has no command-line tools



Audio: 3

Chats: 4

Email: 4

Files: 36

Images: 2

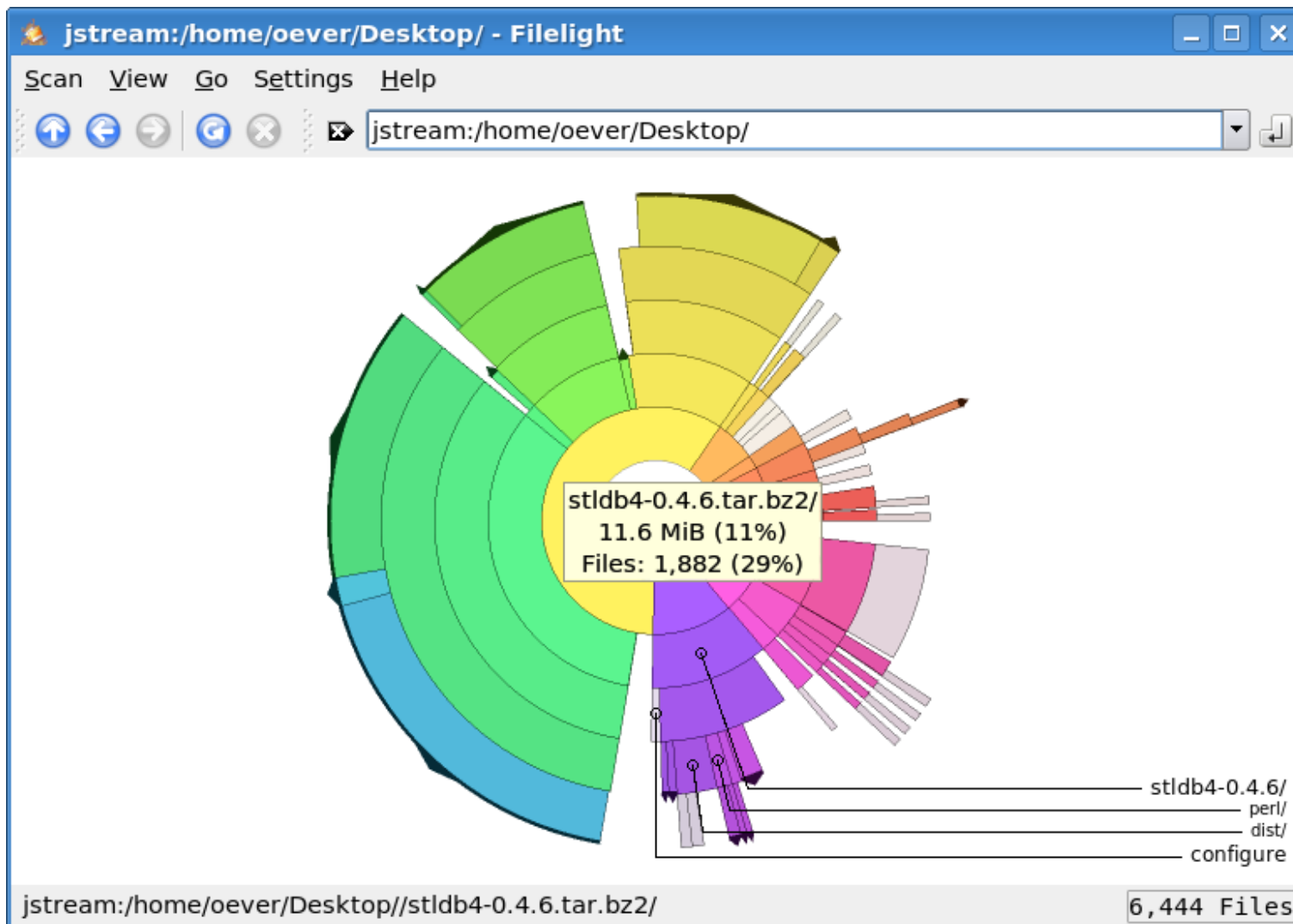
Remote: 2

Source Included: dead link

Video: 3

Web History: 3

Other: 19





Browsing your files



/home/oever/Desktop/ - Filelight [- [□ [×]

Scan View Go Settings Help

↑ ↶ ↷ ↵ ✕ [/home/oever/Desktop/]

stldb4-0.4.6.tar.bz2
2.1 MiB (2%)

/home/oever/Desktop/stldb4-0.4.6.tar.bz2 **3,692 Files**