# Structured Storage and Compound Files

scorpiosoftware.net/2024/11/09/structured-storage-and-compound-files

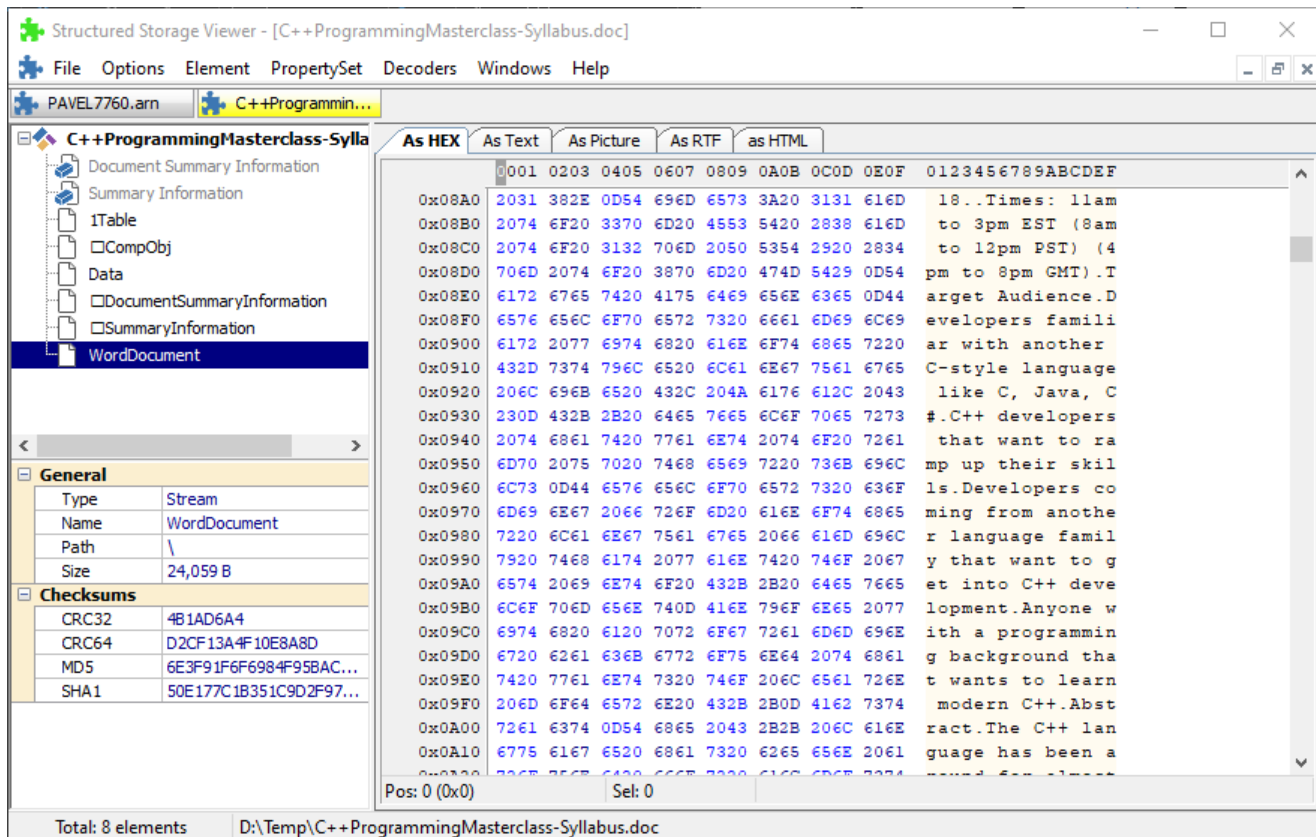View all posts by Pavel Yosifovich                                November 9, 2024

Structured Storage is a Windows technology that abstracts the notions of files and directories behind COM interfaces – mainly `IStorage` and `IStream`. Its primary intent is to provide a way to have a file system hierarchy within a single physical file.
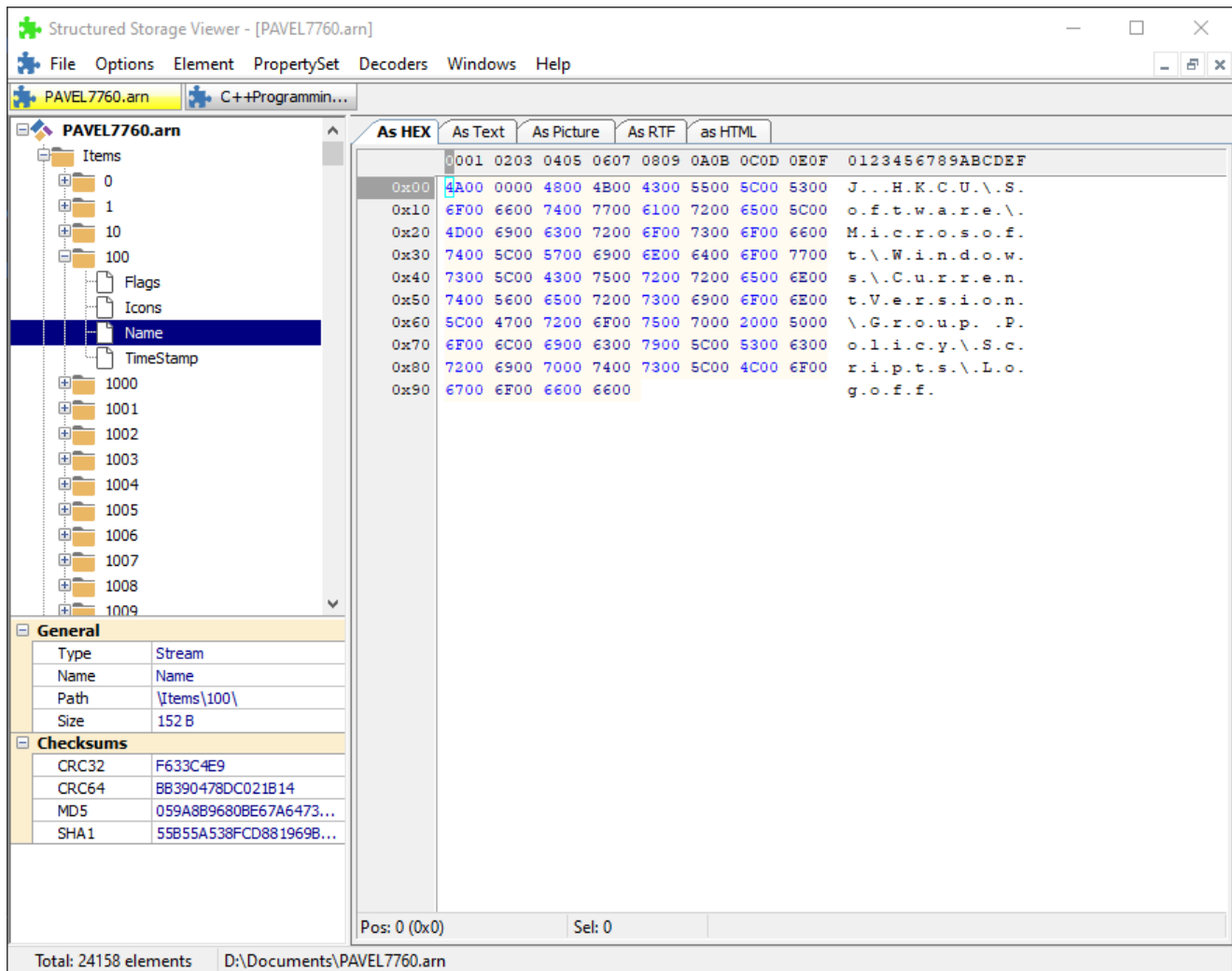
Structured Storage has been around for many years, where its most famous usage was in Microsoft Office files (*.doc, *.ppt, *.xls, etc.) – before Office moved to the extended file formats (*.docx, *.pptx, etc.). Of course, the old formats are still very much supported.

The Structured Storage interfaces (`IStorage` representing a directory, and `IStream` representing a file) are just that – interfaces. To actually use them, some implementation must be available. Windows provided an implementation of Structured Storage called **Compound Files**. These terms are sometime used interchangeably, but the distinction is important: Compound Files is just one implementation of Structured Storage – there could be others. Compound Files does not implement everything that could be implemented based on the defined Structured Storage interfaces, but it implements a lot, definitely enough to make it useful.

You can download an old tool (but still works well) called **SSView**, which can be used to graphically view the contents of physical files that were created by using the Compound File implementation. Here is a screenshot of **SSView**, looking at some DOC file:

Here is a more interesting example – information persisted using Sysinternals *Autoruns* tool (discussed later):

A more interesting hierarchy is clearly visible – although it's all in a single file!

## The Main Interfaces

The `IStorage` interface represents a "directory", that can contain other directories and "files", represented as `IStream` interface implementations. To get started a physical file can be created with `StgCreateStorageEx` or an existing file opened with `StgOpenStorageEx`. Both return an `IStorage` pointer on success. From there, methods on `IStorage` can be called to create or open other directories (storages) and/or files (streams).

The most useful methods on `IStorage` are `CreateStorage`, `CreateStream`, `OpenStorage` and `OpenStream`. Enumeration of storages/streams is possible with `EnumElements`. Here is an example for opening a compound file for read access (filename is from a command line argument):

```
CComPtr<IStorage> spStg;
auto hr = ::StgOpenStorageEx(argv[1], STGM_READ | STGM_SHARE_EXCLUSIVE,
STGFMT_STORAGE, 0, nullptr, nullptr, __uuidof(IStorage),
reinterpret_cast<void**>(&spStg));
if (FAILED(hr)) {
printf("Failed to open file (0x%X)\n", hr);
return hr;
}
```

The following demonstrates enumerating the hierarchy of a given storage, recursively:

```
void EnumItems(IStorage* stg, int indent = 0) {
CComPtr<IEnumSTATSTG> spEnum;
stg->EnumElements(0, nullptr, 0, &spEnum);
if (spEnum == nullptr)
return;
STATSTG stat;
while (S_OK == spEnum->Next(1, &stat, nullptr)) {
if (indent)
printf(std::string(indent, ' ').c_str());
printf("%ws", stat.pwcsName);
if (stat.type == STGTY_STORAGE) {
printf(" [DIR]\n");
CComPtr<IStorage> spSubStg;
stg->OpenStorage(stat.pwcsName, nullptr,
STGM_READ | STGM_SHARE_EXCLUSIVE, 0, 0, &spSubStg);
if (spSubStg)
EnumItems(spSubStg, indent + 1);
}
else
printf(" (%u bytes)\n", stat.cbSize.LowPart);
::CoTaskMemFree(stat.pwcsName);
}
}
```

Each item has a name, but streams ("files") can have data. The `cbSize` member of `STATSTG` returns that size. A stream is just an abstraction over a bunch of bytes. To actually read/write from/to a stream, it needs to be opened with `IStorage::OpenStream` before accessing the data with `IStream::Read`, `IStream::Write` and similar methods.

## More on Streams

The `IStream` interface is used in various places within the Windows API, not just part of Structured Storage. It represents an abstraction over a buffer, that in theory could be anywhere – that's the nice thing about an abstraction. Given an `IStream` pointer, you can read, wrote, seek, copy to another stream, clone, and even commit/revert a transaction, if supported by the implementation. Compound Files, by the way, doesn't support transactions on streams.

Outside of Structured Storage, streams can be obtained in several ways.

The `CreateStreamOnHGlobal` API creates a memory buffer over an optional `HGLOBAL` (can be `NULL` to allocate a new one) and returns an `IStream` pointer to that memory buffer. This is useful when dealing with the clipboard for example, as it requires an `HGLOBAL`, which may not be convenient to work with. By getting an `IStream` pointer, the code can work with it (maybe reading it from another stream, or manually populating it with data), and then calling `GetHGlobalFromStream` to get the underlying `HGLOBAL` before passing it to the clipboard (e.g. `SetClipboardData`).

A stream can also be obtained for a file directly by calling `SHCreateStreamOnFile`, providing a convenient access to file data, abstracted as `IStream`.

Another case where `IStream` makes an appearance is in ActiveX controls persistence.

Yet another example of using `IStream` is as a way to "package" information for COM object's state that would allow creating a proxy to that object from a different apartment by calling `CoMarshalInterThreadInterfaceInStream` (probably the longest-name COM API), that captures the state required (as an `IStream`) to pass to another apartment, where the corresponding `CoGetInterfaceAndReleaseStream` can be called to generate a proxy to the original object, if needed.

## Case Study: Autoruns

Back in 2021 when I was working for the *Sysinternals* team, one of my tasks was to modernize Autoruns from a GUI perspective. I thought I would take this opportunity to do a significant rewrite, so that it would be easier to maintain the tool and improve it as needed. Fortunately, Mark Russinovich was onboard with that, although my time estimate for this project was way off 🙂 But I digress.

One of the features of Autoruns is the ability to save the information provided by the tool so it can be loaded later, possibly on a different machine. This is non-trivial, as some of the information is not easy to persist, such as icons. I don't recall if the old Autoruns persisted them, but I definitely wanted to do so.

The old Autoruns format was sequential in nature, storing structures of data linearly in the file. Any new properties that needed to be added would require offset changes, which forced changing the format "version" and make the correct decisions when reading a file in an various "old" formats.

I wanted to make persistence more flexible, so I decided to change the format completely to be a compound file. With this scheme, adding new properties would not cause any issues – a new stream may be added, but other streams are not disturbed. The code could just ignore properties (could be storages and/or streams) it wasn't aware of. This made the format extensible by definition, immune to any offset changes, and very easy to view using tools, like *SSView*. The above screenshot is from an Autoruns-persisted file.

Persisting icons, by the way, becomes pretty easy, because **ImageList** objects, used by Autoruns to hold collection of icons can be persisted to a stream with a single function call: `ImageList_Write`; very convenient!

## Conclusion

The Structured Storage idea is a powerful one, and the Compound File implementation provided by Windows is pretty good and flexible. One of the reasons Microsoft moved Office to a new format was the need to make files smaller, so the new extended formats are ZIP compressed. Their internal format changed as well, and is not using Compound Files for the most part. A Structured Storage file could be compressed, saving disk space, while still maintaining convenient access using storages and streams.