# Mining the Depths of Excel

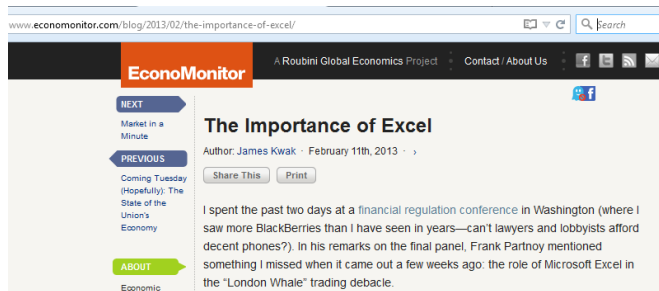## Case Study in Objects and Arrays

# Why Excel?

James Kwak:
(prominent financial blogger?)
*"...Microsoft Excel is one of the greatest, most powerful, most important software applications of all time..."*



**The Importance of Excel**
Author: James Kwak · February 11th, 2013 · ›

I spent the past two days at a financial regulation conference in Washington (where I saw more BlackBerries than I have seen in years—can't lawyers and lobbyists afford decent phones?). In his remarks on the final panel, Frank Partnoy mentioned something I missed when it came out a few weeks ago: the role of Microsoft Excel in the "London Whale" trading debacle.

*"...Excel is everywhere you look in the business world—especially in areas where people are adding up numbers a lot, ... I have a probably untestable hypothesis that, were you to come up with some measure of units of software output,* ***Excel would be the most-used program in the business world****..."*
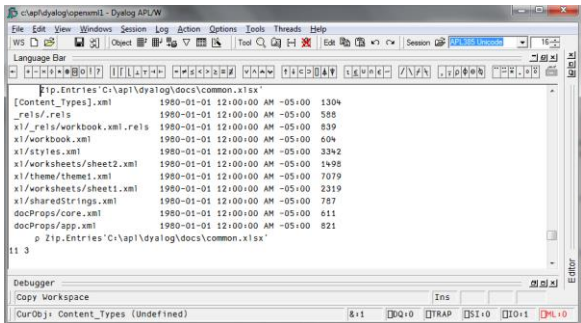
(http://www.economonitor.com/blog/2013/02/the-importance-of-excel/)

# Excel at BCA Research...

- **Lists - keeping track of..., eg. user profiles, data retrieval codes, publication files, etc.**
- **Data sources - (Bloomberg, ThomsonReuters...) make data available as .xlsx or .csv**
- **Interfaces for data collection - downloads and analytical tools are driven by Excel Add-ins**
- **Charts - if no other way to produce**
- **Statistics - if no better way to calculate**
- **Reports - presentation of analyses; lists of things that need attention, etc.**
- **Process control - "table-driven" tasks - determine what to do based on worksheet contents**
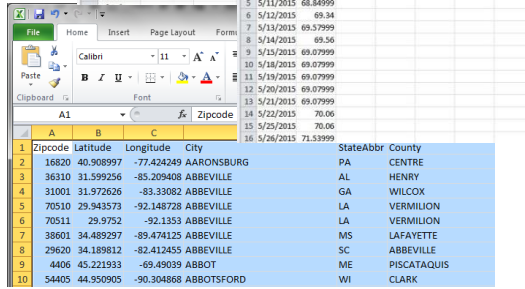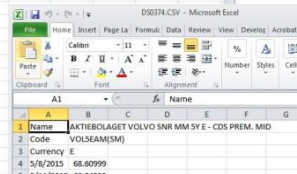
**In short, Excel is extremely important, and EVERYWHERE**

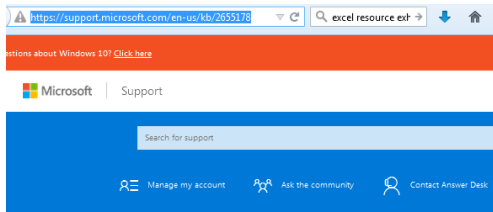# Integration with APL



XLSX
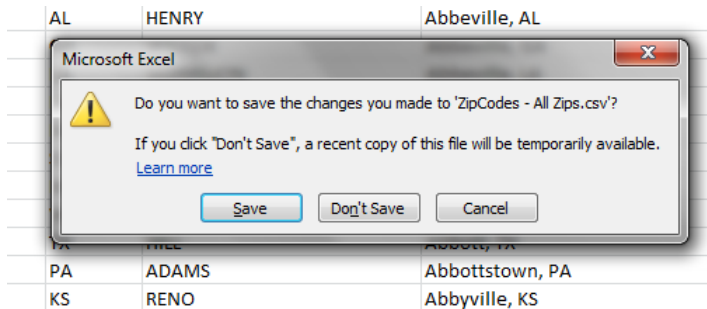CSV
XLS

1000's

NEED for SPEED?
Automation?

# Glitches & Gotchas

- **Installation / Incompatibility issues - APL/Excel**

- **Resource exhaustion - Excel**

- **Automation failures, eg. pop-ups**



Excel cannot complete this task with available resources error, Excel 2010

# Excel and XML

**Office Open XML (OOXML)**
**en.wikipedia.org/wiki/Office_Open_XML**

- **around 2000 - MS began to develop xml-based format for Office documents**

- **2000-2006 - standardization process (ECMA/ISO/IEC)**

- **MS Office 2007 - MS adopts the format as default**

- **\*.xlsx, not \*.xls**

**Office Open XML (also informally known as OOXML or OpenXML) is a <u>zipped, XML-based</u> file format developed by Microsoft for representing spreadsheets, charts, presentations and word processing documents. The format was initially standardised by Ecma (as ECMA-376) and, in later versions, by ISO and IEC (as ISO/IEC 29500). Starting with Microsoft Office 2007...**

### Office Open XML

From Wikipedia, the free encyclopedia

*Not to be confused with Open Office XML*

Article  Talk

WIKIPEDIA
The Free Encyclopedia

Main page

# Support for Office Open XML (non-MS)

**https://en.wikipedia.org/wiki/List_of_software_that_supports_Office_Open_XML**
and
**en.wikipedia.org/wiki/Comparison_of_Office_Open_XML_software#Spreadsheet_documents**

- **Many direct products**

- **other (C#/VB) programming components**

- **Office Open XML standard (ECMA-376 and ISO/IEC 29500:2008)**

Create account    Log in

Article | Talk        Read  Edit  More ▼   Search

**WIKIPEDIA**
The Free Encyclopedia

# List of software that supports Office Open XML

From Wikipedia, the free encyclopedia

Main page
Contents
Featured content
Current events
Random article
Donate to Wikipedia
Wikipedia store

interaction

This is an overview of software support for the Office Open XML

**Office Open XML**
- Office Open XML file formats
- Open Packaging Conventions
- Open Specification Promise
- Office Open XML software
- Comparison of Office Open XML software

# Why not COM? (OLE, ActiveX...)

(http://ramblingcookiemonster.github.io/PSExcel-Intro/)

**"Chances are you have worked with Excel through COM....this isn't supported if you want to use it in an automated solution... "**

*...from Microsoft: (https://support.microsoft.com/en-us/kb/257757)*
*All current versions of Microsoft Office were designed, tested, and configured to run as end-user products on a client workstation. They assume an interactive desktop and user profile. They do not provide the level of reentrancy or security that is necessary to meet the needs of server-side components that are designed to run unattended.*

"Microsoft does not currently recommend, and does not support, Automation of Microsoft Office applications from any unattended, non-interactive client application or component (including ASP, ASP.NET, DCOM, and NT Services), because Office may exhibit unstable behavior and/or deadlock when Office is run in this environment."

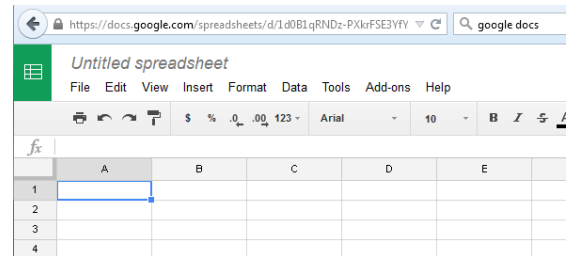**So #1 = Automation;  Plus - above-mentioned issues...**

# Options beyond COM/OLE

## Application?, or API?

- **Who does the UI?**
- **Is there a UI?, etc.**
- **(there's also ODBC?; ADO.Net?)**

## Applications

- **www.openoffice.org - "Calc", formerly by Sun Microsystems, now Apache**
- **Google Docs; Android Apps - work but are features limited?**
- **(MS-OneDrive = Excel Online?)**

# API / Utilities

## Syncfusion XlsIO

### Introduction to Essential XlsIO

**Essential XlsIO** is a 100% native .NET library that generates fully functional **Microsoft Excel Spreadsheets** in **native Excel** format without depending on Microsoft Excel. **Essential XlsIO** is a perfect solution for those who need to read and write Microsoft Excel files from code. It does not require MS Excel to be installed in the report generation machine or server.

**Essential XlsIO** library can be used in any .NET environment including C#, VB.NET, and managed C++. It is a Non-UI component that can be used in **ASP.NET, Windows Forms, WPF, Silverlight, ASP.NET MVC, WinRT, Windows Phone 8** applications, without any change in the API.

# sfExcel

- **excellent toolkit for Dyalog APL (no need to study XlsIO)**
- **Namespace**
- **speed issues if (1000's < 1↑ρMAT)**
- **"DataTable" feature for speedups**
- **some formatting options**
- **data/text only (not "inclusions")**
- **Pierre Gilbert**





**Reading data from the spreadsheet to APL:**

```
xl.GetValue 'A2'        ⍝ Get the valu
xl.GetValue 2 1         ⍝ Get the valu
xl.GetValue 'D5:F8'     ⍝ Get the valu
xl.GetValue 5 4 8 6     ⍝ Get the valu
```

# But what does the C# / VB / VisualStudio crowd do?

## MS Open XML SDK
https://msdn.microsoft.com/en-us/library/office/bb448854.aspx

# Open XML SDK - Download

**API and "Tool"**

Open XML SDK 2.5 for Microsoft Office

| | | |
|---|---|---|
| Language: | **English** | **Download** |

**Note:**There are multiple files available for this download.Once you click on the "Download" button, you will be prompt files you need.

| | |
|---|---|
| **Version:** | **Date Published:** |
| RTW | 11/20/2012 |
| **File Name:** | **File Size:** |
| OpenXMLSDKV25.msi | 2.5 MB |
| OpenXMLSDKToolV25.msi | 24.9 MB |

# Open XML SDK - Tool

- **drill-down into XML structure**

- **structure, not content?**

**file:///C:/Program Files (x86)/Open XML SDK/V2.5/tool/HelpPage.htm**
*"...The **tool** integrates the following functionalities: Automatically generate Open XML SDK code based on document content. User will be able to directly copy, compile and run the code to re-generate the same document or specific parts of the documents..."*
**(? not tried yet)**

# Open XML - Structure



from: http://www.developerfusion.com/article/6170/read-and-write-open-xml-files-ms-office-2007/
- for a good explanation of XML internal structures and specifics, tag definitions, etc.

See also: "Open XML Explained e-book"
http://openxmldeveloper.org/blog/b/openxmldeveloper/archive/2007/08/13/1970.aspx

# Open XML - Structure (from APL)

**:USING System.IO.Compression**  (zip.dyalog by DanB)



```
      Zip.Entries'C:\apl\dyalog\docs\common.xlsx'
[Content_Types].xml           1980-01-01 12:00:00 AM -05:00    1304
_rels/.rels                   1980-01-01 12:00:00 AM -05:00    588
xl/_rels/workbook.xml.rels    1980-01-01 12:00:00 AM -05:00    839
xl/workbook.xml               1980-01-01 12:00:00 AM -05:00    604
xl/styles.xml                 1980-01-01 12:00:00 AM -05:00    3342
xl/worksheets/sheet2.xml      1980-01-01 12:00:00 AM -05:00    1498
xl/theme/theme1.xml           1980-01-01 12:00:00 AM -05:00    7079
xl/worksheets/sheet1.xml      1980-01-01 12:00:00 AM -05:00    2319
xl/sharedStrings.xml          1980-01-01 12:00:00 AM -05:00    787
docProps/core.xml             1980-01-01 12:00:00 AM -05:00    611
docProps/app.xml              1980-01-01 12:00:00 AM -05:00    821
```

# Open XML SDK - .Net Library



**☐USING,←⊂',C:\Program Files\Open XMLSDK\V2.5\lib\DocumentFormat.OpenXml.dll'**

# Open XML SDK (or similar) - Product Integration

**numerous products built with this, or in similar fashion to it by making use of the Open Office Xml format, some free, some not, some out-of-date, eg.**

- **EPPLus - epplus.codeplex.com - uses System.IO.Packaging namespace**
- **github.com/dfinke/ImportExcel - uses EPPLus.dll**
- **GemBoxSoftware - www.gemboxsoftware.com/support/articles/read-write-excel-spreadsheet-net**
- **EASYXLS - www.easyxls.com/**
- **Simple OOXML - http://simpleooxml.codeplex.com/ (2010)**
- **ExcelPackage - excelpackage.codeplex.com - (last update 2007?)**

- **MatLab - http://www.mathworks.com/help/matlab/ref/xlsread.html**
- **R - XLConnect - http://blog.datacamp.com/r-tutorial-read-excel-into-r/**

- **popular Blog:**
  **http://openxmldeveloper.org/**



Home   Blog   Resources   Forums   About Open XML

Welcome to the new
**OpenXMLDeveloper.org**

If you're interested in learning more about Open XML Formats, you're in the right place. OpenXmlDeveloper.org is your source for Open XML Formats news and information, a place to get your questions answered and share your insights as well. It's all free, and open to everyone. This site is regularly updated by a variety of contributors, so check back soon to see what's new!

# Issues (XlsIO and Open XML SDK)

**1. Speed - for sfExcel (XlsIO) - using a "DataTable" may help, BUT...** what about datatype?



sfExcel - GetValue on UsedRange



sfExcel - GetNumber2 (DataTable)

**DataSet & DataTable** (...represents one table of in-memory relational data):
https://msdn.microsoft.com/en-us/library/t31h6yhs%28v=vs.110%29.aspx

# XlsIO - DataTable usage

**http://docs.syncfusion.com/winrt/xlsio/working-with-data#exporting-from-worksheet-to-data-table**

XlsIO / Concepts and Features /
Working with Data

## Exporting from Worksheet to Data Table

It is easy to export the sheet
data to a data table by using the
**ExportDataTable** method of
IWorksheet. This method
allows selection of various data
table options such as include
column names, export formula
calculated values, styles, and
types through the
**ExcelExportDataTableOptions**
enumeration. It has the
following values.

*ExcelExportDataTableOptions enumeration values*

| Members | Description |
|---|---|
| None | Indicates default export to datatable. |
| ColumnNames | Indicates to export ColumnNames to datatable. |
| ComputedFormulaValues | Indicates to export ComputedFormulaValues to datatable. |
| DetectColumnTypes | Indicates that **XlsIO** should try to detect column types. |
| DefaultStyleColumnTypes | When **DetectColumnTypes** is set and this flag is set too, it means that the default column style must be used to detect style. If this flag is not set, but **DetectColumnTypes** is set, then the first cell in the column is used to detect column type. |

# DataTable I-Beams - 2010⍳, 2011⍳

**Dyalog Technical Note #1:** (circa 2010)

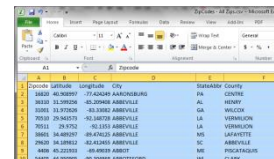### SetDT←2010⌶　　**Set DataTable contents**

### GetDT←2011⌶　　**Get DataTable contents**

"Optimisations to Support the Microsoft.Net System.Data.DataTable Class"

*System.Data.DataTable is a central object in the ADO.NET library in the Microsoft.Net Framework. It is used by DataSet and DataView objects as a container for data which has been extracted from, or will be written to, an ADO.NET data source.*

*"...because we need to "loop" on each row – doing a noticeable amount of work each time. The new I-Beam does all the looping in compiled code..."*

- Rapid conversion of source data to APL arrays
- compromises on datatype, easily overcome in APL
- sfExcel fns: GetNumber2, GetText2, SetText2...
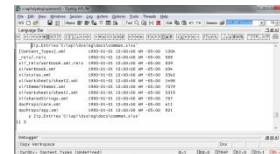- dt.Select(query) = select rows "on the way in" (to APL)?



Excel

DataTable

APL

# Open XML SDK - Issues

- **not many, really, other than complexity, since it offer lots of options**

- **(*guessing...*) similar performance issues as above for large tables** (*when used in "classic OO" manner*)

- **finding online examples is easy (C#), but adapting such code is challenging (for me)**

- **for really large files, blogs indicate workarounds using "Streaming" methods (?)**

- **is it really necessary? ie. if we are just pulling out entire worksheet xml vectors...**

- **may be helpful on inclusions (charts, etc.) - to be investigated; or just sorting out what the workbook contains if unknown?**

- **work in progress**

# **Using Open XML with APL** *- what is the best strategy?*

since we "have" the data...

```
        ρxml ← {Excel/OpenXML worksheet object}...Worksheet.OuterXml
159478
        2000↑xml
```

```
    2000↑3⍴,xml
<x:worksheet xmlns:r="http://schemas.openxmlformats.org/officeDocument/2006/relationships" xmlns:x14="http://schemas.micro
    soft.com/office/spreadsheetml/2009/9/main" xmlns:mc="http://schemas.openxmlformats.org/markup-compatibility/2006" xm
    lns:x="http://schemas.openxmlformats.org/spreadsheetml/2006/main"><x:sheetPr /><x:dimension ref="A1:K352" /><x:sheet
    Views><x:sheetView view="normal" workbookViewId="0"><x:selection pane="topLeft" activeCell="A1" sqref="A1" /></x:she
    etView></x:sheetViews><x:sheetFormatPr baseColWidth="0" defaultRowHeight="15" /><x:sheetData><x:row r="1" spans="1:1
    1"><x:c r="A1" t="s"><x:v>0</x:v></x:c><x:c r="B1" t="s"><x:v>1</x:v></x:c><x:c r="C1" t="s"><x:v>2</x:v></x:c><x:c
    r="D1" t="s"><x:v>3</x:v></x:c><x:c r="E1" t="s"><x:v>4</x:v></x:c><x:c r="F1" t="s"><x:v>5</x:v></x:c><x:c r="G1" t
    ="s"><x:v>6</x:v></x:c><x:c r="H1" t="s"><x:v>7</x:v></x:c><x:c r="I1" t="s"><x:v>8</x:v></x:c><x:c r="J1" t="s"><x:c
    v>9</x:v></x:c><x:c r="K1" t="s"><x:v>10</x:v></x:c></x:row><x:row r="2" spans="1:11"><x:c r="A2" t="s"><x:v>11</x:v
    ></x:c><x:c r="B2" t="s"><x:v>12</x:v></x:c><x:c r="C2" t="s"><x:v>13</x:v></x:c><x:c r="D2"><x:v>19560922</x:v></x:
    c><x:c r="E2" t="s"><x:v>14</x:v></x:c><x:c r="F2" t="s"><x:v>15</x:v></x:c><x:c r="G2"><x:v>42</x:v></x:c><x:c r="H
    2"><x:v>425682</x:v></x:c><x:c r="I2"><x:v>425</x:v></x:c><x:c r="J2"><x:v>77082</x:v></x:c><x:c r="K2" t="s"><x:v>1
    6</x:v></x:c></x:row><x:row r="3" spans="1:11"><x:c r="A3" t="s"><x:v>17</x:v></x:c><x:c r="B3" t="s"><x:v>18</x:v><
    /x:c><x:c r="C3" t="s"><x:v>13</x:v></x:c><x:c r="D3"><x:v>19820324</x:v></x:c><x:c r="E3" t="s"><x:v>19</x:v></x:c>
    <x:c r="F3" t="s"><x:v>20</x:v></x:c><x:c r="G3"><x:v>37</x:v></x:c><x:c r="H3"><x:v>51787</x:v></x:c><x:c r="I3"><x
    :v>384</x:v></x:c><x:c r="J3"><x:v>27292</x:v></x:c><x:c r="K3" t="s"><x:v>21</x:v></x:c></x:row><x:row r="4" spans=
    "1:11"><x:c r="A4" t="s"><x:v>22</x:v></x:c><x:c r="B4" t="s"><x:v>23</x:v></x:c><x:c r="C4" t="s"><x:v>24</x:v></x:
    c><x:c r="D4"><x:v>195
```

*(thanks DanB...)*

# Open XML → APL

**extract items from the "XML vector"...  insert "shared strings"...**

The *row* element defines a new row. Normally the first row in the *sheetData* is the first row in the visible sheet. There are optimizations that can alter this behavior, which is discussed later on in the chapter.  Inside the row you create new cells using the *c* element.

Values for cells can be provided by storing a *v* element inside the cell. Usually the *v* element will contain the current value of the worksheet cell. To store a numeric value in the spreadsheet, all you have to do is to include its value in the *v* element.

```
<worksheet xmlns="http://.../spreadsheetml/2006/main" >
  <sheetData>
    <row>
      <c>
        <v>42</v>
      </c>
    </row>
  </sheetData>
</worksheet>
```
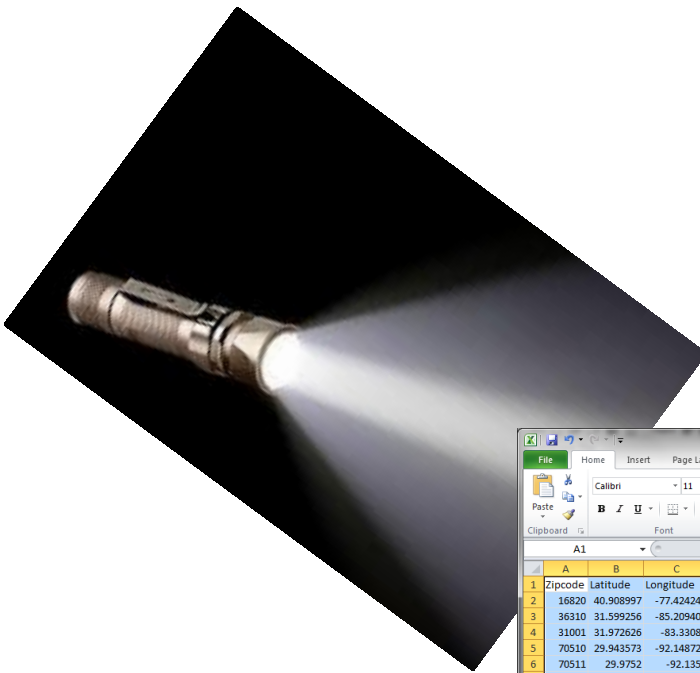
|   | A |
|---|---|
| 1 | 42 |
| 2 |   |

(from: "Open XML Explained":  http://openxmldeveloper.org/..../Open-XML-Explained.pdf)

# Open XML → APL - issues and further work

- ⎕XML vs. direct string manipulations?  regular expressions?
- &amp;  etc. ie. xml escapes
- missing cell references, special-case dates, N/As, etc.
- inclusions  (charts, shapes, tables, etc.)
- selection on (named?) ranges, worksheets, etc.
- clever coding for DataTable and datatype issues (several options)
- APL write to Excel - more to be investigated

**Thank you!**