

Welcome to this live webinar on The large object handling capabilities of SIARD with the Database Preservation Toolkit

Start 10:00

15 June 2023

Audience notes for the Live Webinar



Your cameras have been turned off and microphones muted.



If you have any technical issues during the event, please use the chat function.



Please use the Q&A for questions to speakers. These will be addressed at the end of the event.



Please note that this webinar will be recorded. No attendee personal information will be captured in these recordings.



10:00 – 10:05 **eArchiving Initiative welcome** Jaime Kaminski – eArchiving Initiative training activity lead

10:05 – 10:50 **The large object handling capabilities of SIARD with the Database Preservation Toolkit** István Alföldi – Poliphon

10:50 – 11:00 **Q&A**





Large object handling capabilities of SIARD with Database Preservation Toolkit

István Alföldi, Poliphon (alfi@poliphon.hu)

eArchiving Initiative Training Webinar

Content



- Introduction to DB archiving and SIARD (non-technical)
- Pilot at DH-Lab, Budapest (non-technical)
- LOB handling with SIARD (more technical)



Introduction to DB archiving and SIARD





Database archiving

Other DB elements

- Schemas
- Views
- Keys

• Etc.

- Triggers
- Constraints
- Users, groups, rights



- Start Date
- End Date

♦ ...



SIARD

- Software Independent Archival of Relational Databases (SIARD)
- Open file format for the long-term archiving of relational databases
- Capable of archiving practically all types of DB elements
- Text data based on XML, packaged in a ZIP container files
- Based on standards including Unicode, XML, SQL:2008, URI, ZIP
- Originally developed by the Swiss Federal Archives, later versions in cooperation with the E-ARK projects.
- Now at version 2.2







SIARD

3.1 Use of standards

To ensure that the contents of a database remain interpretable over a long period, the SIARD format is essentially based on two ISO standards: XML and SQL:2008.

ID	Description of requirement	M/O
G_3.1-1	All database content is stored in a collection of files in XML 1.0 format ⁴ that conform to schema definitions according to XML Schema 1.0 ⁵ . The schema definitions and SQL code must in each case conform to SQL:2008 in accordance with ISO/IEC 9075.	М
	The only exceptions are BLOB and CLOB data (Binary Large Objects and Character Large Objects) which are stored in separate binary and text files but are referenced in the XML files.	

3.2 Databases as documents

A relational database is treated as a single document to be archived, so that the references between the data in individual tables are preserved.

ID	Description of requirement	M/O
G_3.2-1	A relational database is archived in a single SIARD file. This file may reference externally stored Large Objects that belong to the database in a larger sense. In rare cases a SIARD file may need to be segmented due to size.	м

4.1 Construction of the SIARD archive file

The SIARD archive file is realised as a ZIP archive.

ID	Description of requirement	M/O
G_4.1-1	The SIARD file is stored as a single, ZIP archive in accordance with the specification published by the company PkWare, version 6.3.2 ⁸ .	М



SIARD – Header

Metadata.xml





SIARD – Mapping







SIARD – Mapping



SIARD specification





SIARD – Content

Table1.xml





SIARD



4 Introduction to Database Archival

- 1. Make sure you know which parts of the database need to be archived. If needed, get in touch with the responsible personnel (e.g. someone from the archive responsible for appraisal).
- 2. Prepare the database for archival: create a new user on the database system that only has read permissions to objects that need to be archived. If needed, create a copy of your database (or certain tables/parts of it) or create views. The database may not be changed during the archival process; otherwise, extraction with SIARD Suite will fail. Never archive from a live system.
- 3. Download the database using SIARD Suite.
- 4. Quality control: check the SIARD file to make sure that everything needed is included, spot check some entries to ensure everything went well.
- 5. Advanced quality control: load SIARD file into a database system again. Run some defined queries on the original database as well as on the archived one and compare results.
- 6. Supplement SIARD file with metadata.
- 7. Define which external documentation needs to be archived together with the SIARD file to ensure comprehensibility of the data (e.g. code tables, system documentation, Entity-Relationship-Diagram, ...).

(Manual SIARD-Suite 2.2, Swiss Federal Archives)



Pilot at DH-Lab, Budapest



DH-Laboratory, ELTE

- National Laboratory for Digital Heritage (DH-Lab)
- Born Digital project
- Archiving research data





- Research database archiving pilot
- Institute of Literary Studies
- DARIAH 2023 Annual Event, at Budapest
- White paper (coming soon)



Introduction

Workflow steps

The National Laboratory for Digital Heritage (DH-Lab) was established in Hungary through an inter-institutional collaboration in September 2020, DH-Lab plans to achieve a set of digital archiving pilot scenarios, starting with a database archiving pilot. This paper provides a summary of the 2022 pilot activities and results. Relational databases are frequently used tools in cultural heritage research. In the scope of the pilot, three research databases created by the institute for Literary Studies of the Hungarian Academy of Sciences have been archived using the information package specifications and database archiving tools of the European E-ARK programme.

Main goals of the pilot: 1. To examine the possibility of preserving research databases in order to create bes practice guides and easy-to-use services in the future.

eArchiving initiative of the European Co

capabilities of SIARD.

Access process



The three digital heritage research databases have been preserved with three different large object (LOB) handling scenarios using the specifications and tools of the eArchiving initiative of the European Commission.



DAME IN F. REAL CO.







eArchiving The E-ARK programme of the Europe

Commission started with the E-ARH project in 2014. The project aimed o synthesise best-practices and develop a core set of interoperabilit specifications for archival operation After the successful conclusion of the E-ARK project, the European Commiss included its outcomes into the Connectin Europe Facility (CEF) programme (2018 2021). Currently, the programme is par of the eArchiving CSP. The E-ARK specifications are based on the Open Archival System (DAIS) reference model. The reference mode sets the digital archiving basics by defining the information packages and the functional units handling and processing these packages. The OAIS model howeve does not define the exact structure of the SIP/AIP/DIP packages. One of the main objectives of E-ARK/eArchiving projects has been the creation and maintenance of information package specificat



Relational Databases (SIARD) is an open file format for the long-term preservation of relational databases. SIARD is based on standards including the ISO standards Unicode, XML, and SQL-2008, the URI Internet standard, and the industry standard ZIP. In the pilot, we examined the large object (LOB) handling capabilities of SIARD Large objects (like images, video etc.) are binary data in a database. Since version 2.0, the SIARD specification supports three possible archival scena for storing LOBs: inline, where LOBs are contained in the DB records, internal, where LOBs are stored outside the DB but within the SIARD file, and external, when LOBs are referenced from SIARD but tored as files outside the SIARD file



POLIPHON



Literature in Western Hungary in 17–19th century

- Novels in Hungary (1730–1836)
- Popular prints in 17th–19th century Hungary
- Literary and scholarly correspondence





Large Objects in the DB

- Images
- Texts
- Formatted HTML pages





Data preparations

- Active research databases \rightarrow Snapshot
- \rightarrow Original SIP
- Formatted HTML pages \rightarrow Searchable text
- PHP user surface \rightarrow Images as context
- Considered the SIP for the Pilot







The Workflow







The Workflow





Used eArchiving components



- SIARD 2.1
- E-ARK SIP/DIP
- CITS SIARD
- Database Preservation Toolkit
- RODA-in
- RODA Repository
- IP Viewer





Large Object (LOB) handling with SIARD



SIARD version history





SIARD v1.0	Original version by the Swiss Federal Archives. Replaced by version 2.1.1: Version 1.0 has been replaced by a more recent and current version. Its use remains possible, but it is recommended to use the present version.
SIARD v2.0	Abrogated: Version 2.0 has been displaced because of errors and ambiguities that might have led to practical problems in the long term. This old issue must no longer be used, and either the present version or version 1.0 must be used.
SIARD v2.1	Replaced by version 2.1.1: Version 2.1 corrects errors and ambiguities in version 2.1. It has been developed by the eCH Fachgruppe on Digital Preservation but is no official standards by the eCH.
SIARD v2.1.1	Replaced by version 2.2: Version 2.1.1 documents the current state of the SIARD file format. It is identical to version 2.1 save for a few precisions in the wording.
SIARD v2.2	Version 2.2 adds support for files outside the database according to part 9 of SQL:2008 (ISO/IEC 9075-9:2008 – SQL/MED) as well as scalability supporting large objects stored outside the SIARD file. Apart from this it is identical to version 2.1.1. It has been developed by the DILCIS Board during the E-ARK3 project.





LOBs in SIARD

Large objects in general are binary data in a database. Binary data is mostly referred to as a binary large object (BLOB) and large character-based data are named character large objects (CLOB).

From version 2.0 SIARD specification supports three possible archival scenarios for storing LOBs:

- 1.) Inline LOBs are contained in the DB records.
- 2.) Internal LOBs are stored outside the DB records but within the SIARD file.
- **3.)** External LOBs are referenced from SIARD but stored as files outside the SIARD file.









DH-Lab archiving scenarios



XVIII-XIX. Century Literature	Data	Images
Popular publications	SIARD	 Inline In the DB records as LOB fields
Novels	SIARD	InternalOutside the DB as files
Author correspondence	SIARD	ExternalOutside the DB as files





DBPTK – Create SIARD file

Create	Open	Manage
This option allows you to create a SIARD file from a supported DBMS.	This option allows you to open a SIARD file.	This option allows you to open, edit, validate, migrate, or visualize the information about SIARD file previously ingested.
CREATE (2)	OPEN O	MANAGE #

European

Commission

	db dat	abase, toolkit eservation		
	[Home Create Manage Preferences H	elp	
Create	Ope	☆ Home > ※ Create SIARD - Connection		
This option allows you to create a SIARD file	This option allows you t	DBMS	General SSH Tunn	el
from a supported DBMS.		S JDBC	Hostname *	localhost
CREATE Z'	OPEN	Microsoft Access		The name of the database server host (e.g. localhost)
		S Microsoft SQL Server	Port number	3306
		🛢 MySQL		The server port number
		S Oracle	Username *	zkalo
		S PostgreSQL	Password *	i ne name of the user to use in connection
		Progress Openedge	Password	The password of the user to use in connection
		Sybase	Database *	itiponyva
				The name of the database to connect
			🛃 Disable Encrypt	tion
			Use to turn off encry;	ption in the connection
			TEST CONNECTION	
	l			_



DBPTK – Create SIARD file







DBPTK – Create SIARD file



Create	Open	Manage
This option allows you to create a SIARD file from a supported DBMS.	This option allows you to open a SIARD file.	This option allows you to open, edit, validate, migrate, or visualize the information about SIARD file previously ingested.
CREATE Q'	OPEN O	MANAGE #







DBPTK – View SIARD file

Create	Open	Manage
This option allows you to create a SIARD file from a supported DBMS.	This option allows you to open a SIARD file.	This option allows you to open, edit, validate, migrate, or visualize the information about SIARD file previously ingested.
CREATE (2'	OPEN O	MANAGE #

European

Commission

DBPTK Desktop			- 0
me Create Manage Pref itipon (7b038724-eb7c-4671-b2d6-3c58 Ponyva adatbázis - képek a	ferences Help IYVA 843c52d0a) adatbázisban, html, link segédtáblában, inner LOB		
Archival date: Archiver: Archiver contact: Client machine:	2022-04-08 02:00 (UTC+2) zkalo kaloz@poliphon.hu BornDigitalVM	Database product:MySQL 8.0.28-0Data origin timespan:1782-1895Data owner:Kőszeghy PéterProducer Application:DBPTK Develop	ubuntu0.20.04.3 per (version 2.9.10)
SIARD See details about the ar columns and other techn Version: SIARD Path: PONYW Size: 5566.8	chived format and manage the metadata of tables, nical elements of the database. O standard version 2.1 (SEE SPECIFICATION) A.SIARD MB	Validation See details about the validation process. Status: Valid Details: ⊘ 27 successes ② no failures ▶ 42 skippe Warnings: There are no warnings Validated at: 2023-06-12 07:09 (UTC+2) Validate by: DBPTK developer 2.9.10 (SEE DETAILS)	Browsing See details about the database status. Manage the database metadata by configuring table, columns and data. Status: Pre-loaded



DBPTK – Validate SIARD file

Create	Open	Manage
This option allows you to create a SIARD file from a supported DBMS.	This option allows you to open a SIARD file.	This option allows you to open, edit, validate, migrate, or visualize the information about SIARD file previously ingested.
CREATE Q'	OPEN O	MANAGE E

	-				
b DBPTK Deskt	top			-	
ome Create	Manage Preferences Help)			
🖌 Home 🗲	🖺 Databases > 🛢 itipony	va 🕽 🖄 <u>Validation</u>			
🗘 Va	alidation				
/alidates the nformation (e SIARD against its specific needed to understand why	ation. The validator sh the requirement faile	nows information about which the requirements have passed and which one have fa d.	iled. In case of a failed requirement, the report file generated	ontains th
Database I	Name:	itiponyva		SIARD specification:	SIARD-2
Requireme	ents that passed:	27		Additional checks specification:	OPEN
Requireme	ents that failed:	0		Report:	OPEN
Number of	f errors:	0			
Number of	f warnings:	0	Validation		
Number of	f skipped:	42	Validation successfully completed for database "itiponyva"		
Status:		Valid	ci	LOSE	
		_		Scroll to	the end 🤇
T_6.4-2	Validation finish on pa	ath: content/schema1,	/table37/table37.xml	O	
T_6.4-2	Validation running on	path: content/schema	1/table38/table38.xml		
T_6.4-2	Validation finish on pa	ath: content/schema1,	/table38/table38.xml	O	
	Validation running on	path: content/schema	1/table39/table39.xml		
T_6.4-2	-				

DBPTK – Search in SIARD file

DBPTK Desktop								- 0		
ne Create Manage Preferences Help										
Home > 📑 Databases > 🛢 itilevelezes > 0	Q. <u>Search</u>									
, Filter sidebar	0									
		ch all ree	cords							
Information									_	
Q Search all records	Arany								٩	
	🖽 archiv	levelezes f	oto 2022							
archiv_b1_html_2022	id	le	evelezes_id	img_order		img_name	img			
archiv_br_tink_2022	147	8	2	0		files/Arany János Toldy Ferenc	inek, Pe			
archiv levelezes html 2022	148	8	2	1		files/Arany János Toldy Ferencnek, Pe				
archiv levelezes link 2022	149	8	3	0		files/Arany János Toldy Ferencnek, Pe				
archiv namespace link 2022	150	8	4	0		files/Arany János Toldy Ferencnek, Pe				
□	151	8	5	0		files/Ipolyi Arnold Arany Jánosnak, Po				
 ⊞ b1	152	8	5	1		files/Ipolyi Arnold Arany Jánosnak, Pe				
🖽 Ь2	1-6 of 6 🔹 🕨					······································				
🖽 helyseg										
🖽 levelezes										
levelezes_audit	—									
levelezes_uggroups	田 levelez	es								
levelezes_ugmembers	Id	elo ki	Papa ki	bonnan	datum	orzobely	ielzet	leiras		
levelezes_ugrights	02	0.0_1.1		Doct	1965 11 17	MTA Köpuntára ás lafa		Fogalmaruánu	A.c.24	
levelezes_users	02			Pest	1005-11-17	MTA Könyvtára és info		Fa an lana an vénu	A	
levelezés_settings	83			Pest	1805-12-13	MTA Konyvtara es inro	MIA KIK KE RAL 1805/	Fogalmazvany	Ага	
🖽 message	84			Pest	1865-12-13	MTA Konyvtara es Info	MIA KIK KERAL 1865/	<table style="b</td> <td>)OX-</td>)OX-	
🖽 namespace	85			Pest	1865-04-29	MTA Könyvtára és Info	MTA KIK KERAL 1865/	Levél, Ipolyi Ar	nolo	
	1-4 of 4								•	

Q. 🥥

9



n 🌀 🔱 🖾 hun 🌐 d× 🗁

2023. 06. 13.



Large Objects - Internal

		j				Œ) lob5	
🔓 Start Page	Schema Editor	onyva					ponyva.siard + content	t + <u>schema1</u> + <u>table6</u>
C C E itiponyva	Execute	+ i += 1+ i					Table40	record10.bin
Events (0) Functions (4)		SC □				table5		Type: BIN File
 ☐ Links (0) ▶ ☐ Tables (40) ▶ Views (0) 	id ponyvak	id ponyvakataszter_id img_order img_name					lob5	Type: BIN File
	2 695 3 696	54	1 files/IMGP8863 Klio_1vzkivuu.jpg 0 files/B1850_2.02_rvqchqio.jpg	JPEG Image			table7	record41.bin
	4 697 5 698 6 699	81 81 224	1 files/B1850_2.04_8ts4lgjl.jpg 2 files/B1850_2.03_42yzccju.jpg 0 files/IMGP9221 ill_4iug1h0m.jpg	JPEG Image JPEG Image JPEG Image		:	table8	Iype: BIN File
Tree View XSL Output xml table xsi:sche xmlns:xs row cl cl cl cl cl cl cl cl cl cl	emaLocation si text ile ength igest igest	<pre>version= http://w http://w http://w 694 54 0 files/IM content/: 370837 fed0a8024 SHA-256</pre>	"1.0" encoding="UTF-8" ww.admin.ch/xmlns/siard/2 ww.admin.ch/xmlns/siard/2 ww.w3.org/2001/XMLSchema- GP9658 ill_b6jf9qc3.jpg schema1/table6/lob5/recor ea47cbc04a5b4d59b85125014	/schema1/table6.xsd /schema1/table6.xsd instance d1.bin f979dd6134e84fbdeae1	table6.			
i row i row					1	美国的		

💐 | 🗧 | ponyva - WinZip Pro

A CONTRACT OF A CONTRACT.

Zip

Unzip

File

Manage

Backup/Clean

Tools

View



Large Objects – Internal

😭 Start Page 💋 🔎 Scher	ma Editor	sol itiponyva				
🖉 🕒 🔪 🛢 itiponyva 🗘 Exec	ute					
	5 č 🖷		[[] [] [] [), 🖵 🖏		
Events (0)	1 SELE	CT * FROM archiv	/_ponyvakatasz	ter_illusztraciok_2	2022;	
Functions (4)						
	id	ponyvakataszter_id	img_order	img_name		
Views (0)	1 694	4 54	0 files/	MGP9658 ill_b6jf9qc3.jpg	JPEG Image	
	2 695	5 54	1 files/	MGP8863 Klio_1vzkivuu.jpg	JPEG Image	
	3 696	5 81	0 files/	B1850_2.02_rvqchqio.jpg	JPEG Image	
	4 697	81	1 files/	B1850_2.04_8ts4lgjl.jpg	JPEG Image	
	5 698	8 81	2 files/	B1850_2.03_42yzccju.jpg	JPEG Image	
	6 699	9 124	0 files/	MGP9221 ill_4iug1h0m.jpg	JPEG Image	
	7 700	224	1 files/	MGP9222 ill_r8huolwf.jpg	JPEG Image	
	8 701	224	2 files/	MGP9227 ill_0kjjyiix.jpg	JPEG Image	
	9 702	2 225	0 files/	MCP8830 ill_r7h4ndxe ind	IPEG Image	
	10 703	do DBF	PTK Desktor			
	11 704	Home	Create Manage Pr	eferences Help		
		9 Q B	Information Search all records	a a	= rchiv_po illusztrac	nyvakatas iok 2022
		:≡	Tables		earch	_
				raria_ncmi_2022		
			archiv_bibliog	rafia_link_2022 k. 2022	a 🔪	ponyvakataszter_id
				K_2022	594	54
			archiv ponyva	akataszter_timap_2	595	54
			⊞	6	596	81
			archiv_ponyvakat	aszter_illusztraciok	597	81
			archiv_szerzo	_link_2022 6	598	81
			🖽 besorolas	6	599	224
			🖽 bibliografia	-	700	224
			eloford_pony	vaban	704	224
			🖽 incipit_moder	'n '	01	224



– 🗆 🗙

ataszter MANAGE ARCHIV_PONYVAKATASZTER_ILLUSZTRACIOK_2022 OPTIONS D22 advanced V C

 Files/IMGCP9658 ill_b6jF9
 Download

 files/IMGCP8863 Klio_1vzk
 Download

 files/B1850_2.02_rvqchqi
 Download

 files/B1850_2.04_8ts4lgjl
 Download

 files/B1850_2.03_42yzccj
 Download

 files/IMGCP9221 ill_4iug1
 Download

files/IMGP9222 ill_r8huol Download

files/IMGP9227 ill_0kjjyii> Download

European Commission





Large Objects – External



26 22



- External file structure
- SIARD 2.1 does not define the file structure for External LOBs
- There is a recommendation from the E-ARK project: *Recommendation for storing large objects outside the SIARD file (E-ARK project, 2017)*
- SIARD 2.2 defines the structure for External LOBs (but DBPTK has not implemented it yet)



External file structure





DATALINK type

- SIARD 2.2 also supports the DATALINK type
- DATALINK type
 - DATALINK is a special SQL type intended to store URLs in database according to SQL:2008 part 9 SQL/MED (ISO/IEC 9075-9:2008).
 - It contains a reference to a LOB in a system outside of the RDB, but partly controlled by the RDBMS. In SIARD it is treated as a LOB with information about the original path.
- Rarely used due to limited support by major DB vendors in the past.



Keep in touch



ec.europa.eu/



europa.eu/



@EU_Commission



M



@EuropeanCommission



European Commission



europeancommission

@EuropeanCommission





Thank you

Contact



https://e-ark4all.eu/



info@e-ark-foundation.com

@EU_eArchiving



https://www.linkedin.com/company/euearchiving-initiative



https://www.youtube.com/@e-ark