

## Publication

# A Distributed Archival Network for Process-Oriented Autonomic Long-Term Digital Preservation

## Thesis (Dissertationen, Habilitationen)

**ID** 2185193

**Author** Subotic, Ivan

**Author at UniBasel** [Subotic, Ivan](#) ;

**Year** 2013

**Title** A Distributed Archival Network for Process-Oriented Autonomic Long-Term Digital Preservation

**Type of Thesis** Dissertation;

**Start of thesis** 01.02.2008

**End of thesis** 05.04.2013

**Name of University** University of Basel

**Name of Faculty** Philosophisch-Naturwissenschaftliche Fakultät;

**Supervisor(s) / Fachvertreter/in** Schuldt, Heiko ;

The rapidly growing production of digital data, together with their increasing importance and essential demands for their longevity, urgently require systems that provide reliable long-term preservation of digital objects. These systems have to ensure guaranteed availability, integrity, authenticity, and interpretability over the course of the preservation, where the preservation period may last for several years, for instance in business or scientific applications, the lifetime of a human in medical applications, up to potentially unlimited time-spans for preservation in cultural heritage digital libraries. This means that all kinds of technical problems (network, software or hardware failures) need to be reliably handled and that the evolution of data formats is supported. At the same time, systems need to scale with the volume of data to be archived. Thus, long-term digital preservation systems have to be inherently distributed to allow content to be replicated. Institutions with long-term archiving needs for the preservation of digital data, have to collaborate in order to build a highly reliable and available, geographically distributed Internet-based digital archiving system. By employing distributed systems technologies be it for the creation of a small cooperating network of few institutions with limited resources, or a large network with many nodes providing combined potentially vast amounts of globally distributed resources, the challenges lie in the autonomic, efficient, and fault-tolerant use of these resources without a centralized global coordinator.

We present novel concepts for a distributed long-term preservation system for digital data, with a focus on long-term preservation as required by archives, museums, research communities, or the corporate sector. These concepts are the result of combining distributed, autonomic, and process oriented computing, with requirements from the digital preservation community regarding special system, user, and metadata functionality. Originating from this fusion, our novel concepts are the main ingredients of the described system model, consisting of a data model, and different processes. At the data level, support is provided for complex data objects, management of collections, annotations, and arbitrary links between digital objects. At process level, our proposed archiving system model supports automated processes that provide dynamic replication, consistency checks, and automated recovery of the archived digital objects utilizing autonomic behavior governed by preservation policies without any centralized coordinator in a fully distributed network. This allows for an efficient and fault-tolerant use of the resources provided in the network. Further, we present a prototype implementation of the DISTARNET (DISTRIBUTed ARchival NETwork) System, a distributed long-term digital preservation solution, which utilizes the described novel concepts. While implementing the described data model and processes, our implementation is additionally governed by considerations such as fault-tolerance on the node level, maintainability

and extendability, and longterm use of the system, which all flow into the described system architecture, and resulting implementation. Subsequently, we then perform an evaluation of the implemented prototype and the underlying concepts, with the use of realistic scenarios. The evaluation consists of two parts. In the first part, we define and employ a benchmark geared towards triple stores, in which we evaluate the feasibility and the constraints of using triple stores for RDF-based metadata storage and management in the context of long-term preservation systems. In the second part, we perform a qualitative and quantitative evaluation of the DISTARNET system prototype implementation. The former looking at the correct execution of the developed processes, and the later looking at the performance of the system regarding the overall archiving storage capacity and scalability of the system.

**Full Text on edoc ;**