# Challenges of Longterm Preservation of Digital Data

*DAAD / UFPR, Dept. Computer Science*

*Research Seminar*

*15<sup>th</sup> March 2011, Curitiba*

Dr. Dirk von Suchodoletz

Faculty of Engineering, University Freiburg, Germany

UNI
FREIBURG

# Overview of the Things to Come

- My Background

- Threats on digital objects

- Developments in this field

  – International initiatives and Projects

  – Existing components and ideas
- Preservation action and Emulation

  – Formalizing the rendering requirements

  – Software archiving

- Workflow integration and automation

# My Background

- Lecturer/researcher at the professorship of communication systems

  – Lectures, seminars and student projects on computer networking and communication systems

  – PhD on "Longterm preservation of dynamic digital objects" in 2008

  – Participating in large scale EU integration project PLANETS

  – Research at National Archives of New Zealand beginning this year

# My Background

- **Chair in Communication Systems**

  – Rather small entity with a small number of scientific assistants/lecturers working in different fields

  – Offers lectures and seminars in Internet Working, Telecommunication Systems, Network Technology; in cooperation with Max-Planck-Institute for Foreign and international Crimial Law seminars on Internet&Law

  – Got into domain of DP via PhD theses on Emulation and PLANETS project – Preservation and Longterm Access to NETworked Services

# The Professorship

- Digital Preservation involvements

    - EU PLANETS Project (finished)

    - Founding member of the *Open Planets Foundation*

    - Cooperation with the National Library and National Archives in The Netherlands

    - Member of the German *nestor* initiative and founding member of its *Emulation WG*

    - Long cooperation with the *Computer Games Museum* opened last weekend in Berlin

    - Supervision of a number of Bachelor- and Master theses in this field

# Short Risk Analysis of Digital Objects

- Threatened on several layers

  - Physical

  - Technological

  - Intellectual

# Physical Risks

- Decay of media

  - More risky are removable, because often "uncovered", unprotected media
    - Optical media like CDROM, DVD, BlueRay
    - DVD is pretty risky even for brand new pressed disks (bacteria eating up the glue layer)
    - CD for several years, no much known on BlueRay (general not a good idea)

  - Less problematic for hard-drives – but number of mechanical problems here

# Physical Risks

- Deprecation of connectors and standards

  - Anyone knows the MFM/RLL controller?

  - Last disk of this standard might be produced just 20 years ago

  - SCSI – pretty old standard – anyone succeeded in connecting the old SCSI-disk to a modern SCSI-320 controller?

  - Same for IDE $\rightarrow$ SATA

# Technological Risks

- Even if the bit-stream recovery was successful (disk completely copied to a modern, accessible medium)

- Rapid changes of the typical work-desks

  - Changes of (G)UI concepts (IBM 286 and Apple iPad just 25 years apart)

  - Different hardware architectures

  - Different operating systems

  - What to do with the old file formats (WordStar, AmiPro, Lotus-1-2-3, WordPerfect, old MS-Office)?

# Intellectual Risks

- Extremely relevant for archival data

- Changes in contextual knowledge

    – Missing or incomplete documentation

    – Lost context of single objects or groups of objects, linked objects

    – Ambiguous data formats and descriptions

    – Changes in terminology and basic assumptions

# Additional Risks

- Human errors

- Technical, machine breakdowns (e.g. a hot summer maneuvered the Freiburg universities computer center to nearly shutdown of most of the server machines in the "air-conditioned" hall)

- Catastrophes

- Security flaws, forgery of data, sabotage ...

# Additional Risks

- Typical risks in many memory institutions
  - Under-financed IT services
  - Understaffed, ill-qualified personnel
  - Missing rights management, authorization, …
  - Missing Know-How

# From Risks to Solutions

- Research Domain of Digital Preservation pretty young but differentiates into sub domains

- Our specialization is preservation action: Emulation
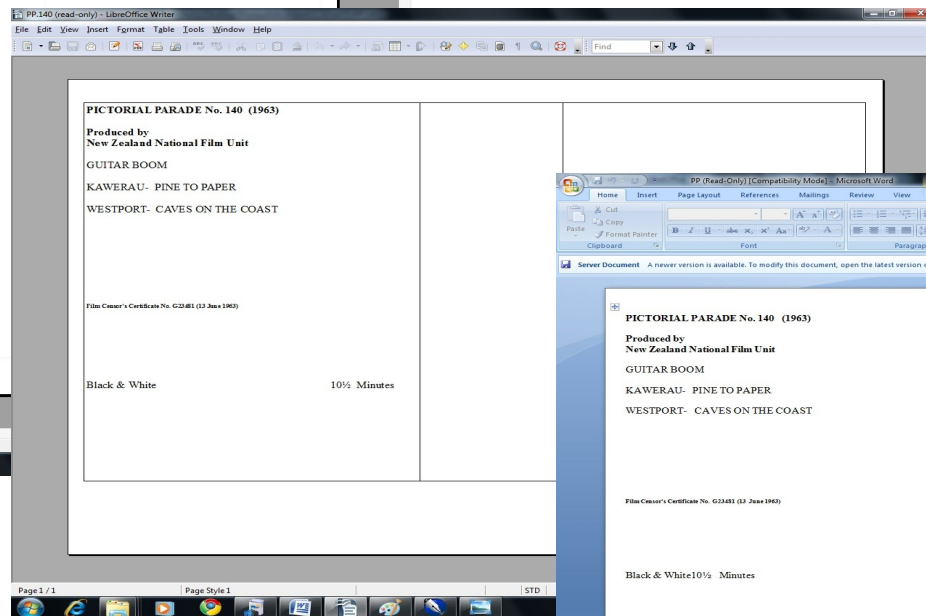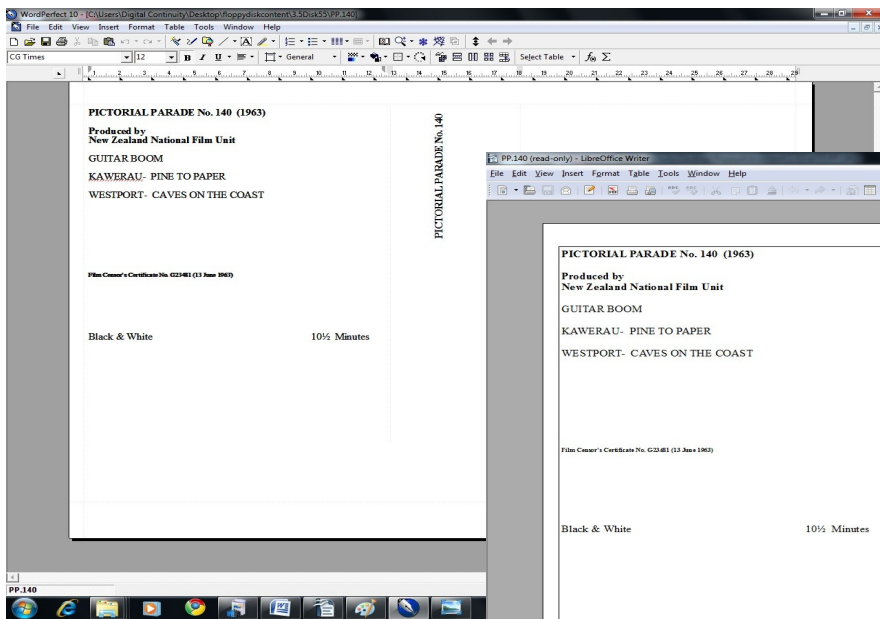
# Dynamic Objects & Authenticity

- Preservation challenges

  - Digital objects require software / hardware environments to be accessed

  - Environments change over the time and obsolete most of digital material

  - Mainline strategy: *Migration*

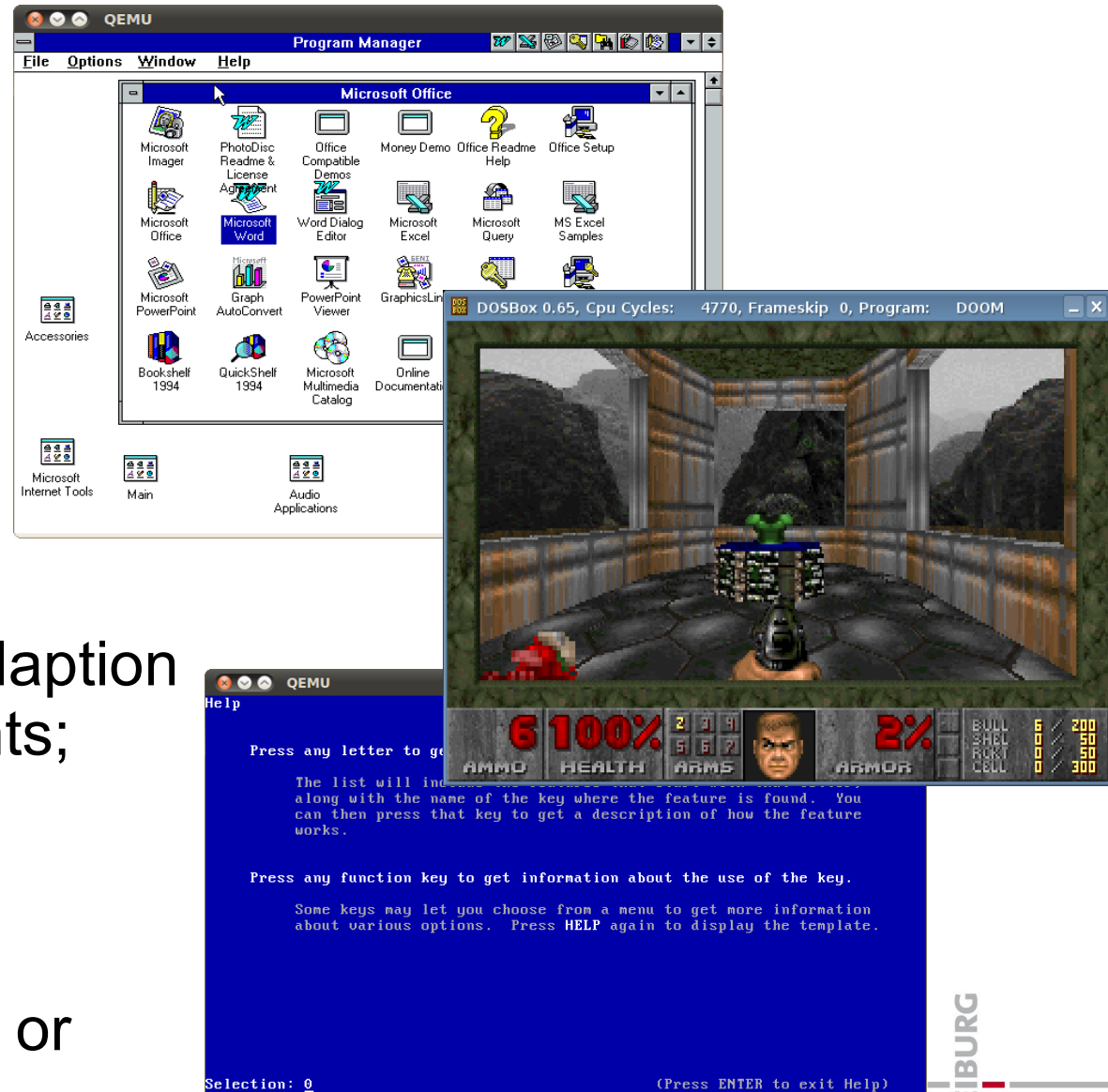  - Risky to rely on it exclusivly

  - Not suitable for all object types

# Authenticity Tests / Experiments

Problems easily spotted: Rendered a test corpus in different original applications/environments

# Dynamic Digital Objects

- **Objects like**
  - Applications
  - Operating systems
  - Databases
- **Non-linear, user inter-action, multiple views**
- **No real option:**
  - Printing of source, adaption to recent environments; even **if** source code available
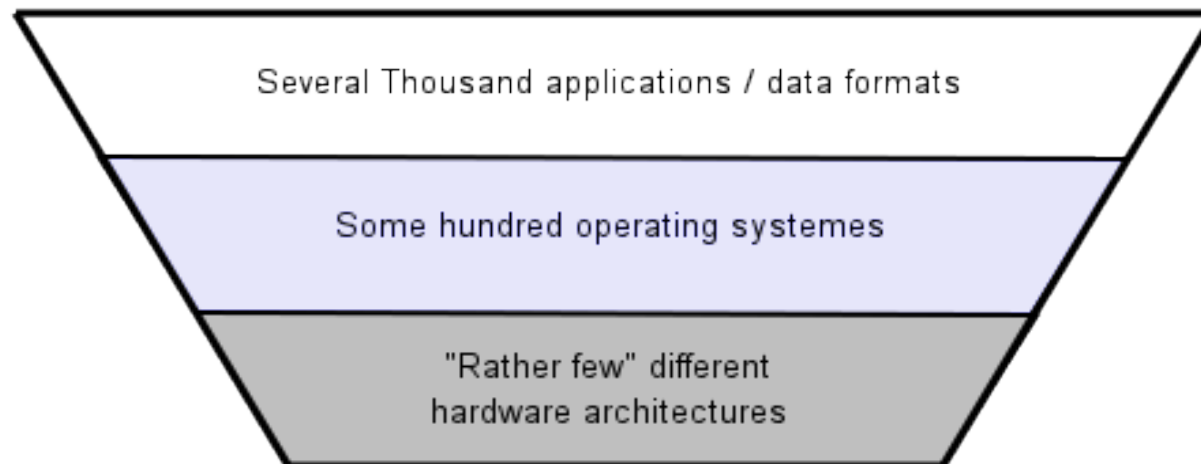  - Video-recording, screenshots of game or application session
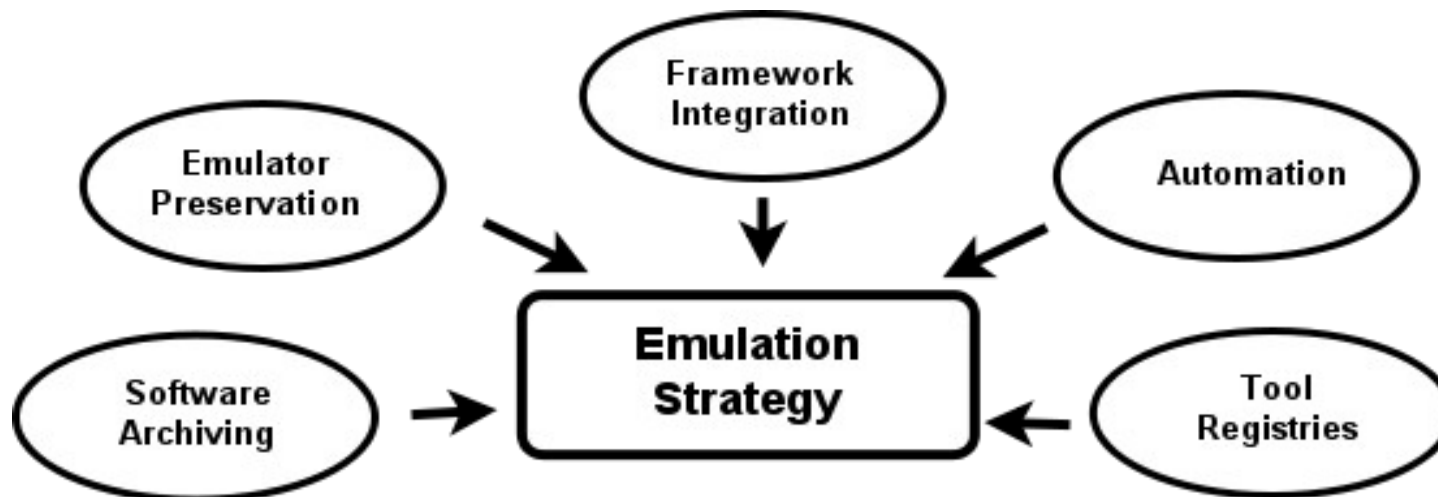
# Different Approach – Emulation

- No changes on object, but reproduction of original environment
  - Emulators around for quite a while, supplemented by virtualization
  - Can operate on different layers of software/hardware stack
  - Number of objects to cover differs significantly; thus hardware layer seems very attractive to focus on

Several Thousand applications / data formats

Some hundred operating systemes

"Rather few" different
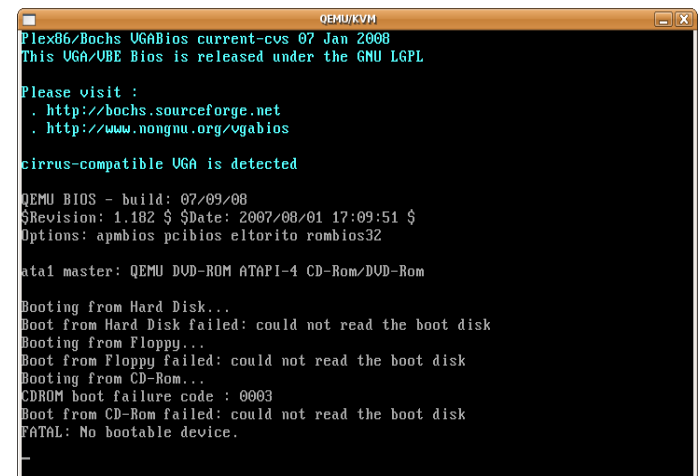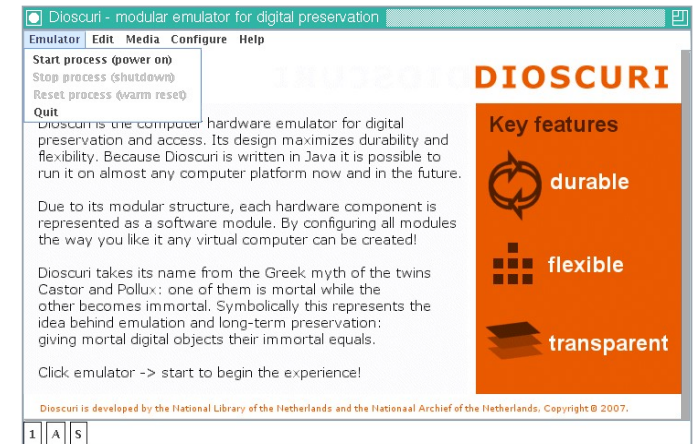hardware architectures

# Emulation Strategy

- More complex approach involving larger number of additional software components, complexities
    - Standalone emulation does not help much
    - Different sectors of ongoing research
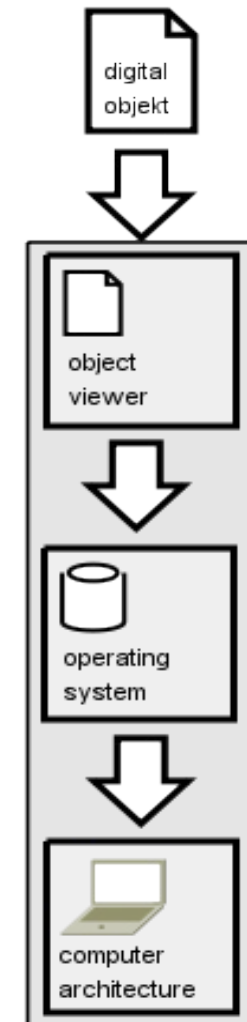    - Number of open issues

# Emulator Examples

- Dioscuri X86 emulator recreating an 286, 386 PC of the early 1990th
  - Java programming language, modular approach – components like disk, floppy, VGA, CPU, RAM put together to form the machine
  - Running DOS and Windows 3.0
- QEMU – using popular C program-ming language multi architecture emulator for X86, PPC, Sparc, ...
- Both Open Source – no vendor lock-in, adaptable
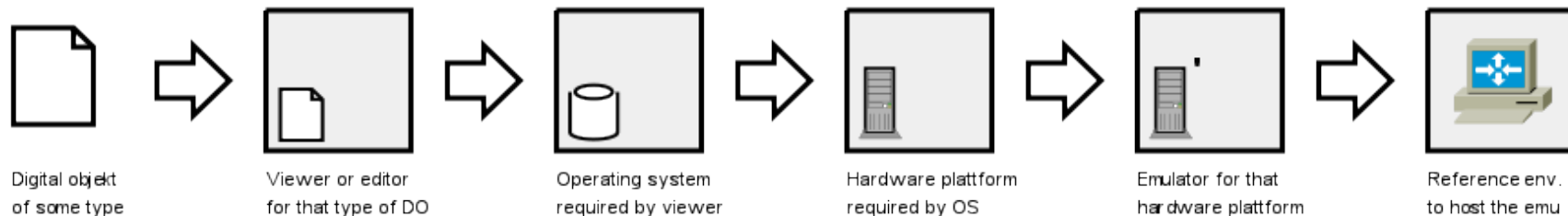
# Requirements for Emulation

- Emulation not working just on its own – additional software is required
- Emulation approach requires recreation of ancient hardware / software environments for access / execution
  - E.g. spreadsheet document requires the proper spreadsheet application for interpretation and displaying
  - Spreadsheet software is dependent on an operation system
  - Operating system was programmed for a very specific or a range of hardware architectures
  - Additional components like fonts might be needed for range of documents, especially for non-latin typesets

digital objekt

object viewer

operating system
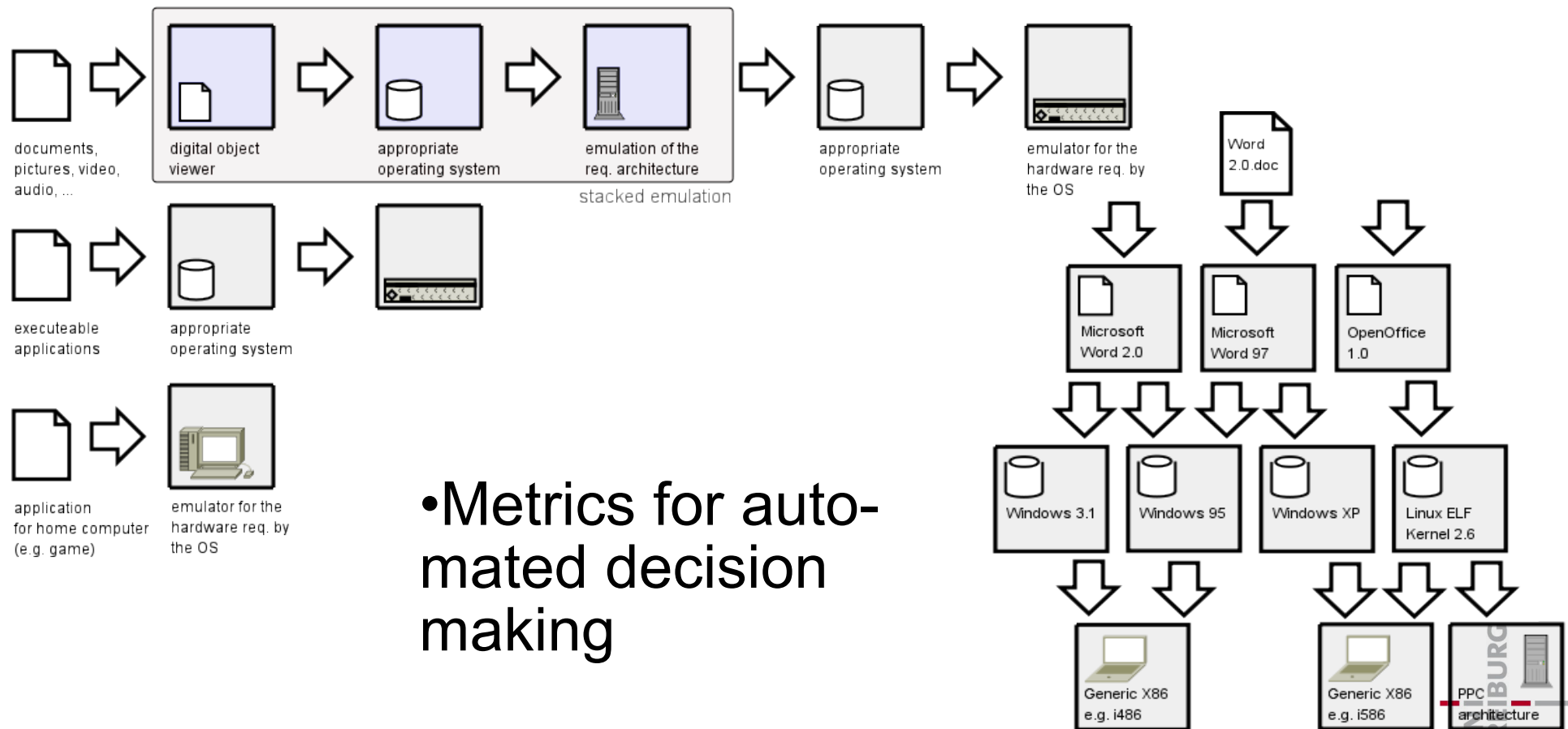
computer architecture

# Formalization of Requirements

- View Path – pathway from object to specific environment

- Formalization needed – view path as the requirements to be followed to actually access, display the object of interest

- Reference environment – specifically defined software hardware combination for object access, rendering

- Concept to describe dependencies between objects

Digital objekt of some type → Viewer or editor for that type of DO → Operating system required by viewer → Hardware plattform required by OS → Emulator for that hardware plattform → Reference env. to host the emu
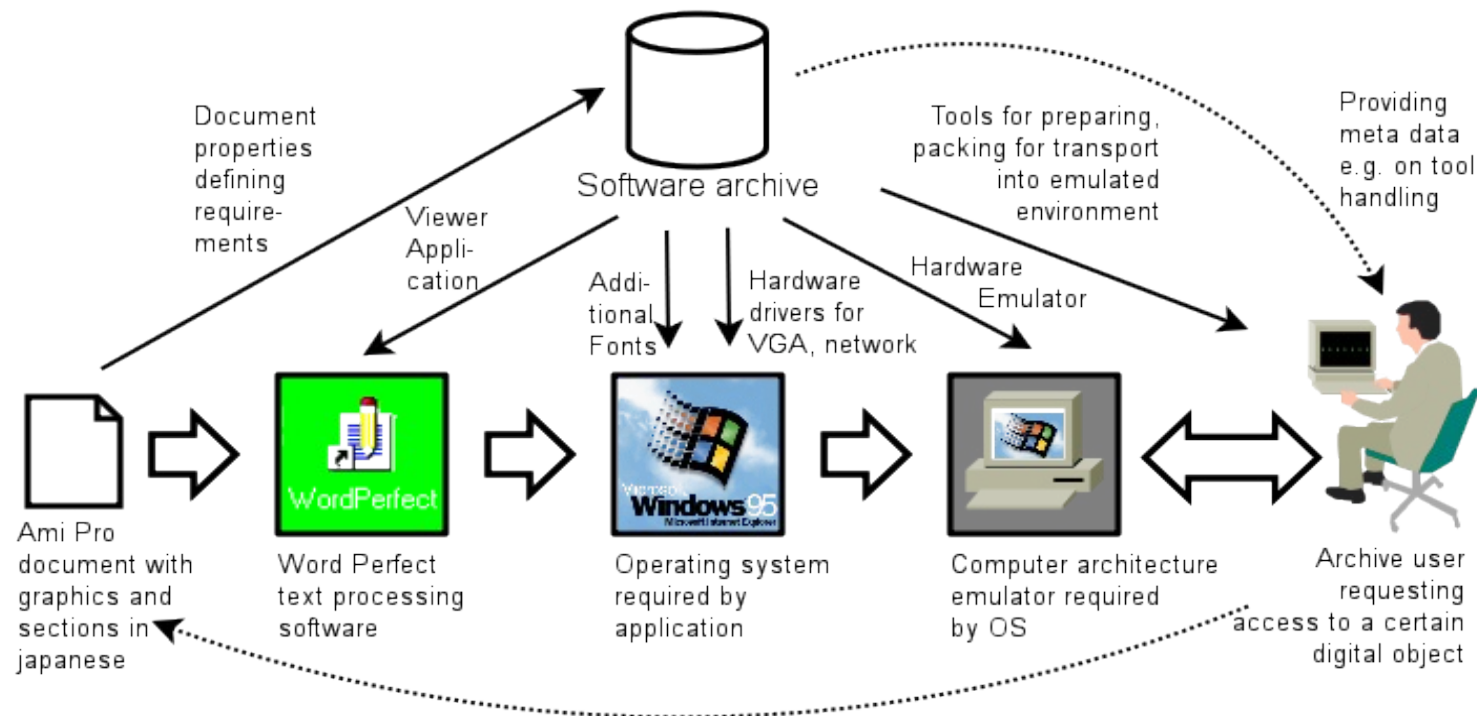
# Formalization of Requirements

• View Path dynamic – depend on regarded object, actual working environment, emulator preservation strategy; often multiple options



• Metrics for auto-mated decision making

# Software Archiving

- Software archive containing all necessary additional single objects or for production systems prefabricated view-paths
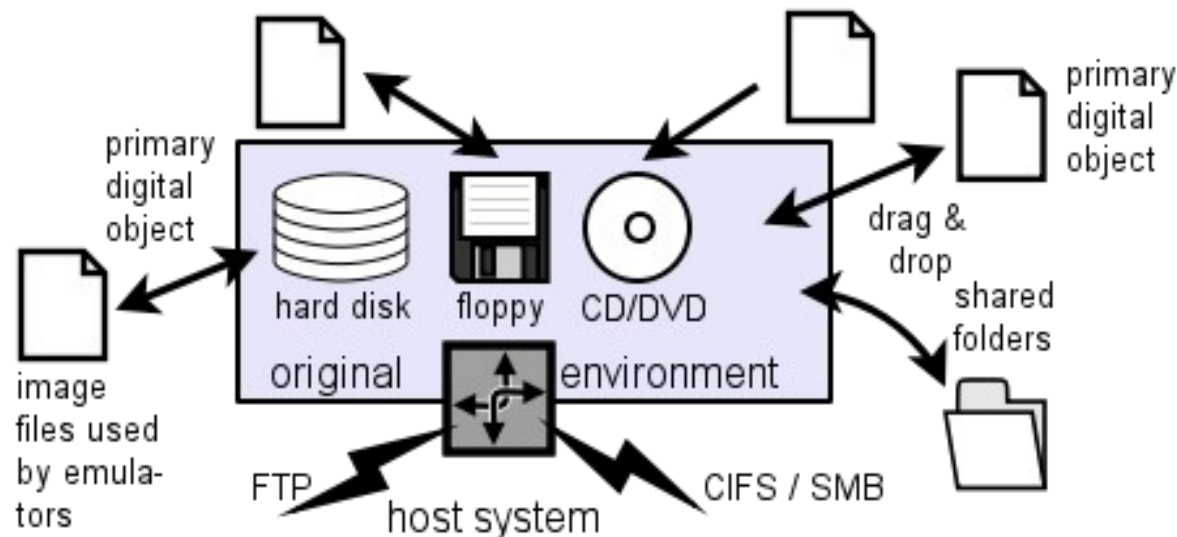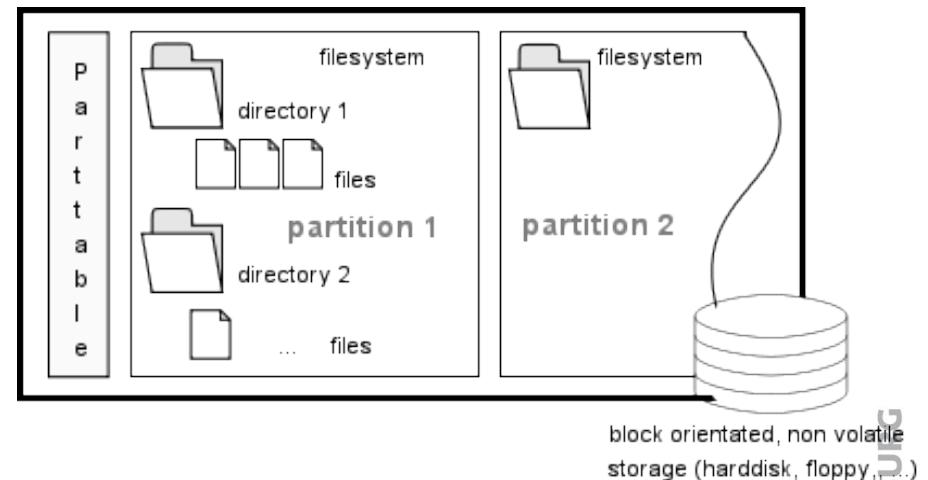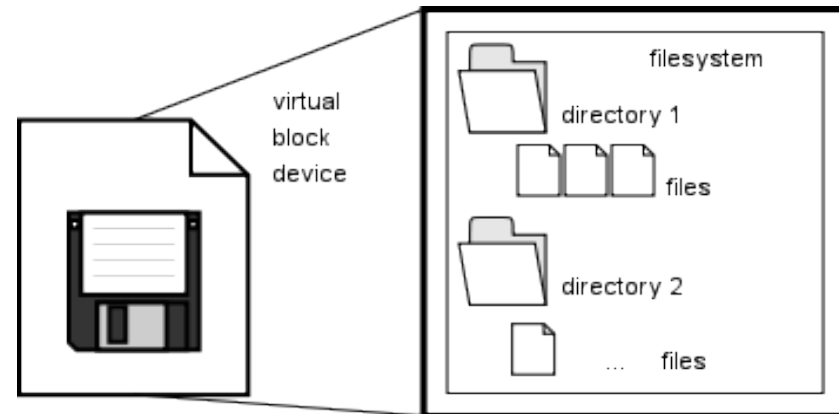
# Additional Information

- Additional information and metadata needed in software archive
  - Application handbooks
  - Howtos and trouble shooting guides
  - Application update packages
  - License keys, access codes
- Depending on object
  - Fonts for documents
  - Codecs for video, audio
  - Software extensions like DirectX, OpenGL libraries

# Data Exchange with Emulators

- Object is to be transported into emulation environment
  - Different ways: After or during enviroment setup



primary digital object

image files used by emulators

hard disk   floppy   CD/DVD

original   environment

FTP   host system   CIFS / SMB

drag & drop

shared folders

primary digital object

- Means of object transport
  - Virtual optical (ISO) or floppy disks as images
  - Disk container files
  - Network connections like FTP, SMB/CIFS
  - "Shared Folders" (as e.g. found in VMware or VirtualBox)
  - Copy&Paste (e.g. text areas in Dioscuri)

# Transport Containers

- Data transport requires fromats understood by the target environment, e.g.
  - Floppy disks, ubiquious in for many platforms for a rather long period
  - Images easy to create and store
  - Optical disks: ISO images well understood by many emus
- More complex
  - Container files of the several emulators
  - Creator tools required
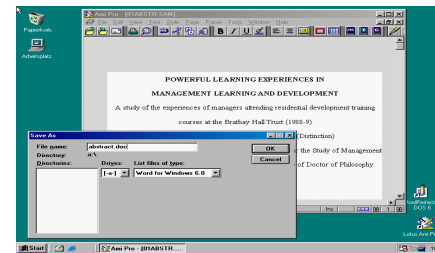  - Adding objects to disk container files before emulators started



filesystem

directory 1

files

directory 2

... files

virtual block device



Partable

filesystem

directory 1

files

directory 2

... files

partition 1

filesystem

partition 2

block orientated, non volatile storage (harddisk, floppy, ...)

# Challenge: Access to Emulation

- Groundworks laid – how to make emulation accessible?

- Emulation environments

    - Often deal with outdated concepts of software interaction

    - Typically complex and require specific knowledge

    - Require depending on the digital object to be rendered or executed a bunch of additional software components which may need prior installation

# Enabling Access to Emulation

- Major goal is to allow non-technical users access to those services an easy to use, abstract interface is required

- During PLANETS project a prototype for emulation wrapping created – GRATE

- Different emulators like Dioscuri, MESS, QEMU, Hatari and others put into a single networked application

# Knowledge and Automation

- GRATE focuses on traditional human interaction model, but

  - Requires certain knowledge getting more and more uncommon for todays users

  - Taking system images of emulated environments for granted

  - Handling only limited, prefabricated VPs

- Unsuitable for integration into non-interactive large scale preservation work-flows

# Knowledge and Automation

- Typical applications most digital objects produced with are interactive

- Standard migration work-flows like opening a document and save it in a different format require a human user to type or point&click

- Such manual procedure sub-optimal for e.g. mass migration scenarios

- Next step: A method to replace the human-interaction in GRATE with generic recording and monitoring

# Automation of Interaction

- Define an interactive work-flow as ordered list of interactive events passed on (e.g. mouse and keyboard events)

- Each event is linked with a precondition and an expected outcome

- Built the solution on top of the VNC-Play tool, which offers visual synchro-nization points

# Framework Integration

- PLANETS – Preservation and Longterm Access to NETworked Services

    – Offers a set of standardized Web services like Characterization, View, Validation, Comparison, Migration, ...

- Defines a set of APIs Web services need to conform to

# Goals

- Emulation services should allow

  - Occasional users to view digital objects and compare digital objects in their original environment

  - Occasional users to experience ancient (graphical) interactive user environments

  - Documentation and preservation of user interactions and interactive processes in ancient user environments

  - Automated migration of files using the original application in emulation

# Required Services

- **After reviewing these goals**
  - **View service to allow traditional interactive access to objects**
  - **Automated migration by emulation service**

# Emulation View Service

- A generic PLANETS view service takes a digital object and returns an URI pointing to the rendered result

- If the digital object requires a running rendering engine the service offers methods for querying and sending

- View service developed allows access to already configured and ready-made software environments

# Emulation View Service

- Implemented as a PLANETS Web service

    - Accepts a list of digital objects

    - Wraps them into a CD image

    - Makes them available for running operating system

- User is able

    - Explore the original environment

    - Use within the original application

    - Allows visual comparison for migrated objects

    - Do manual migrations by saving or printing

- Process can be generalized, recorded

# Emulation View Service

# Migration by Emulation Service

- Good for viewing but not for large scale preservation tasks

- Second important service for PLANETS using emulation

- Interface expects a digital object as input, a designated output format (PUID) and an optional list of service parameters

- Outcome will be a successfully transformed object or an error message

# Actual State of Workflows

- Got some promising results for simple migrations like loading an AMI Pro document and send it to virtual PDF printer

- At the moment pretty expensive regarding time and compute resources:

  - For every object complete cycle from mounting, loading, system execution and shutdown required

  - In future: start system once and handle multiple objects in succession

UNI
FREIBURG

# Black Box Emulation Environment

- Pretty much "black box" at the moment

  - Not yet reliable

  - Migration processes just started difficult to monitor

  - Unknown execution time

  - Unreliable behavior – system might take infinitely or might crash – mostly unobserved

# Measurement and Evaluation

- Very scarce measurement and evaluation options in todays emulators

  - Difficult to calculate runtime, effort, host system resources

  - Tedious to observe file operation e.g. when the object is finished to be processed and completely written to disk again

- Without measurements and metrics for significant characteristics comparisons of workflows and different emulators pretty impossible
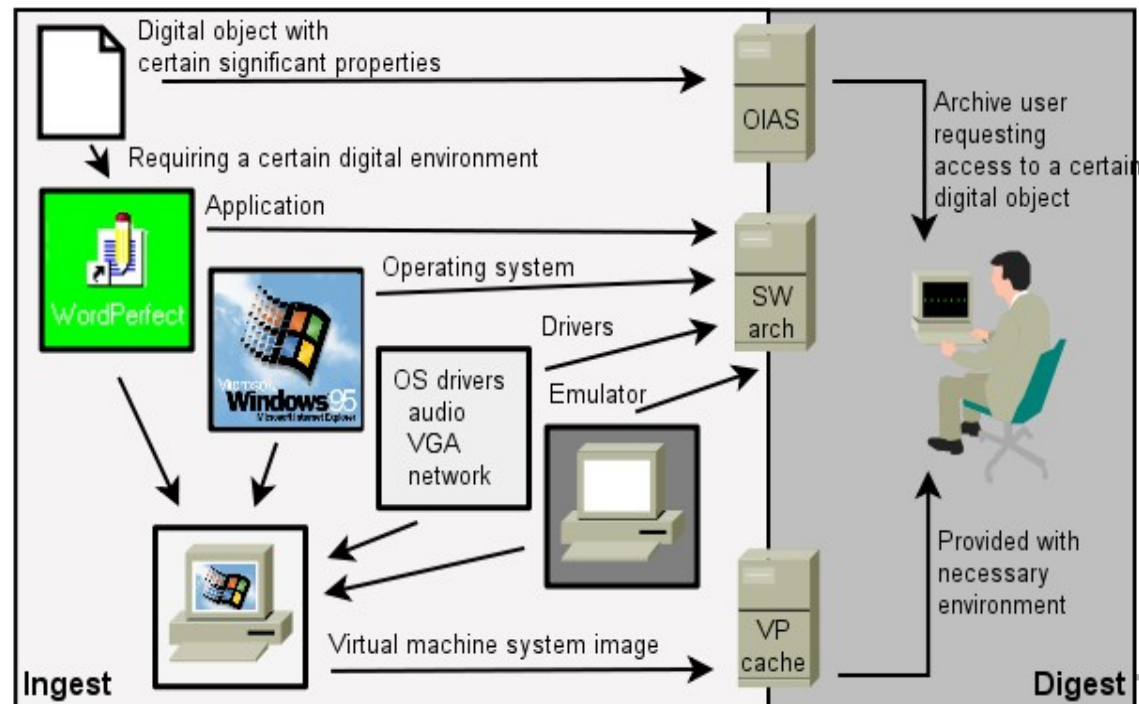
# Major Challenges

- To make migration via emulation workflows comparable

  - Monitoring and evaluation framework needed
  - Have metrics for certain emulation characteristics
  - Test or prove completeness of emulation

# Outlook: General Integration

- Integrate software archive into preservation work-flows

    - Check software list on object ingest

    - Store single software components

    - Documentation

- Preserve knowledge by storing workflow recordings and complete emulation environments

# Ongoing Research / Theses

- Map view path to software archiving / library to preservation workflows

- Large scale workflow automation

- Creating sample reading room workstation for object access through emulation (services)

- Automated emulator testing

- Define future emulator requirements
    - Control APIs (input, automation, monitoring, …)
    - Stable presentation towards original environments

- General: Long-term stable software platforms

UNI
FREIBURG

# Thank you for your Attention!

## Questions / Comments?

**Dirk von Suchodoletz**

Chair on Communication Systems

Faculty of Engineering

University of Freiburg / Germany

*dirk.von.suchodoletz@rz.uni-freiburg.de*