UC SANTA CRUZ

# Deep Store
## Problems and Solutions for the Next Generation of Archival Storage

April 1, 2004

Lawrence You, Kristal Pollack, Svetlana Kagan, Darrell Long

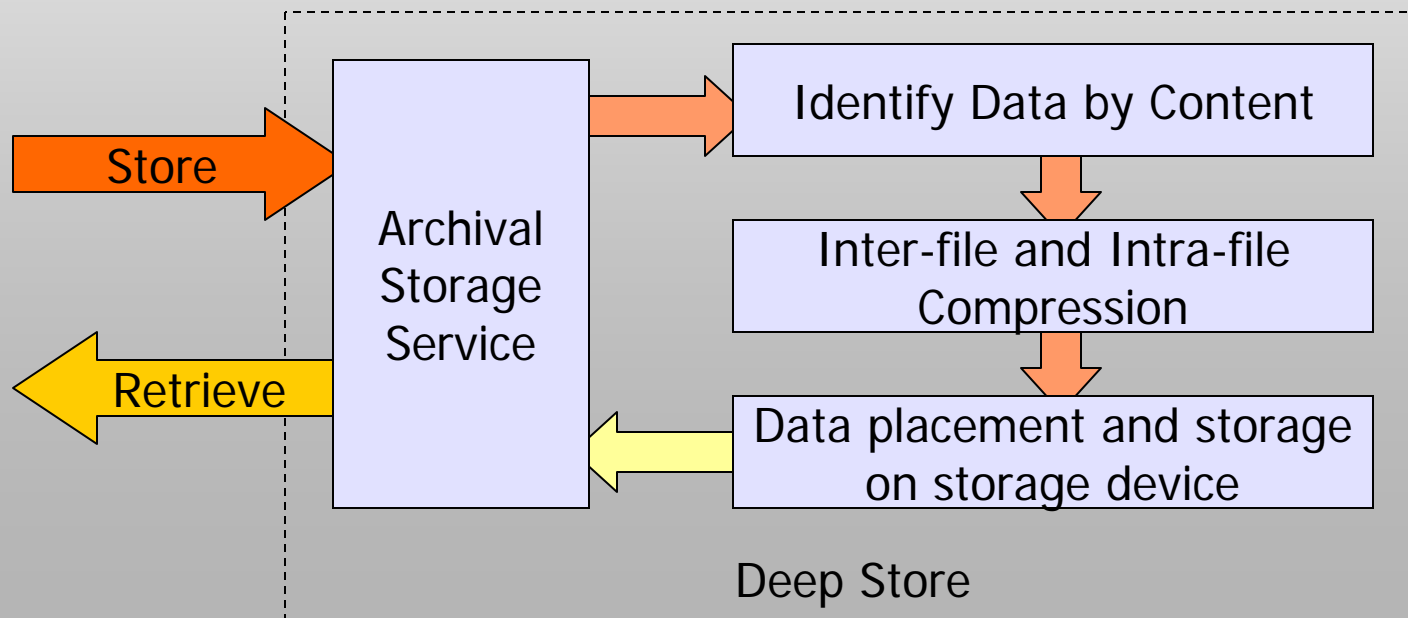University Of California, Santa Cruz

# Problem 1: Cost

**Problems:** Growing volumes of reference (archival) data

Managed disk storage is still more expensive than tape

# Efficient Archival Data Storage

**Solutions:** Improve storage efficiency by eliminating redundancy

Exploit duplication and similarity with inter-file and intra-file compression

# Problem 2: Managing Content

**Problems:** Archival content lives and dies with applications and systems
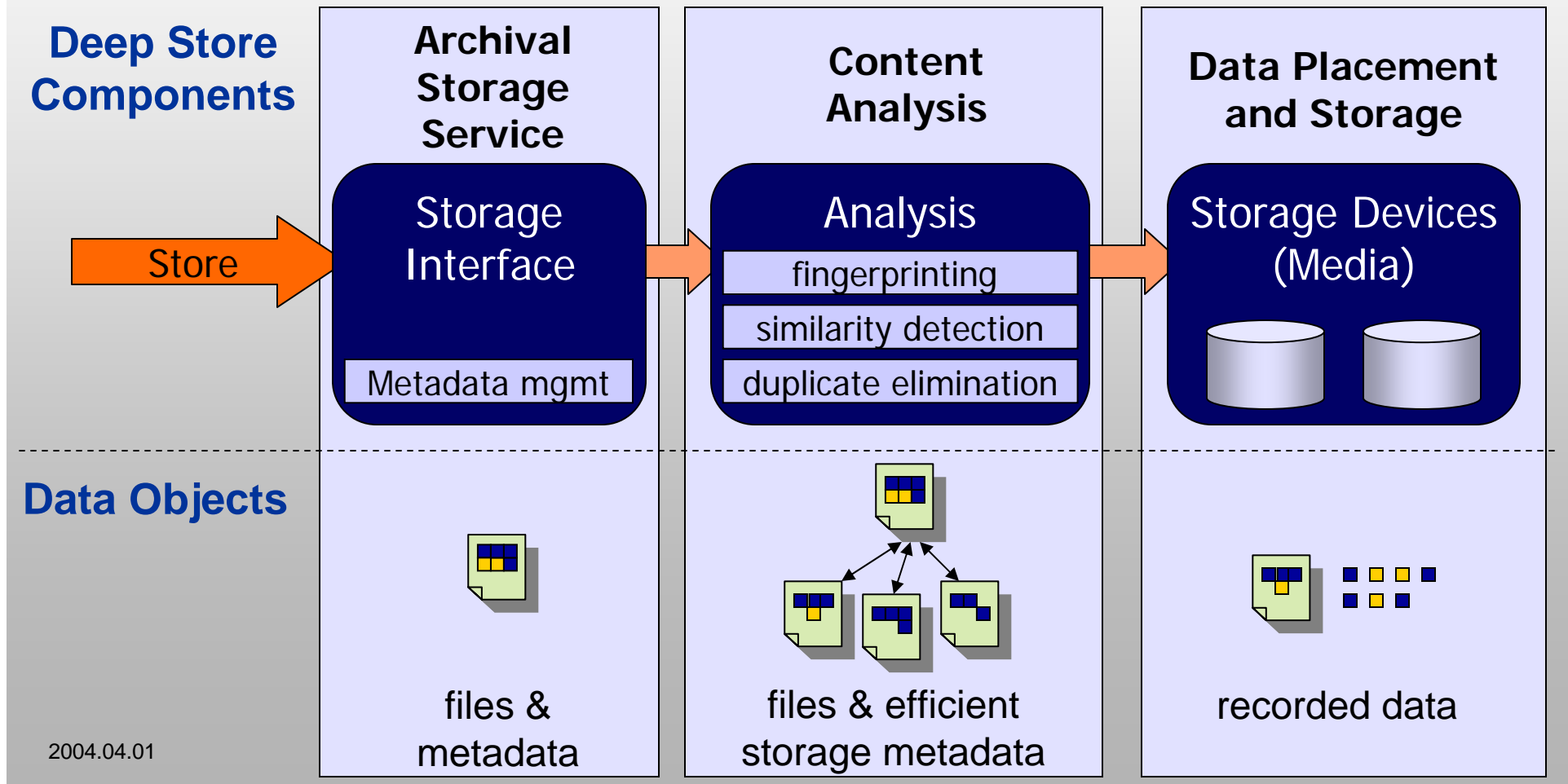
Today's storage systems do little to help future-proof the content

Reference data must live beyond systems and storage devices

# Managing Content

**Solutions:** Manage content with metadata

Create an archival storage interface, replicate, and actively self-monitor

# Problem 3: Performance

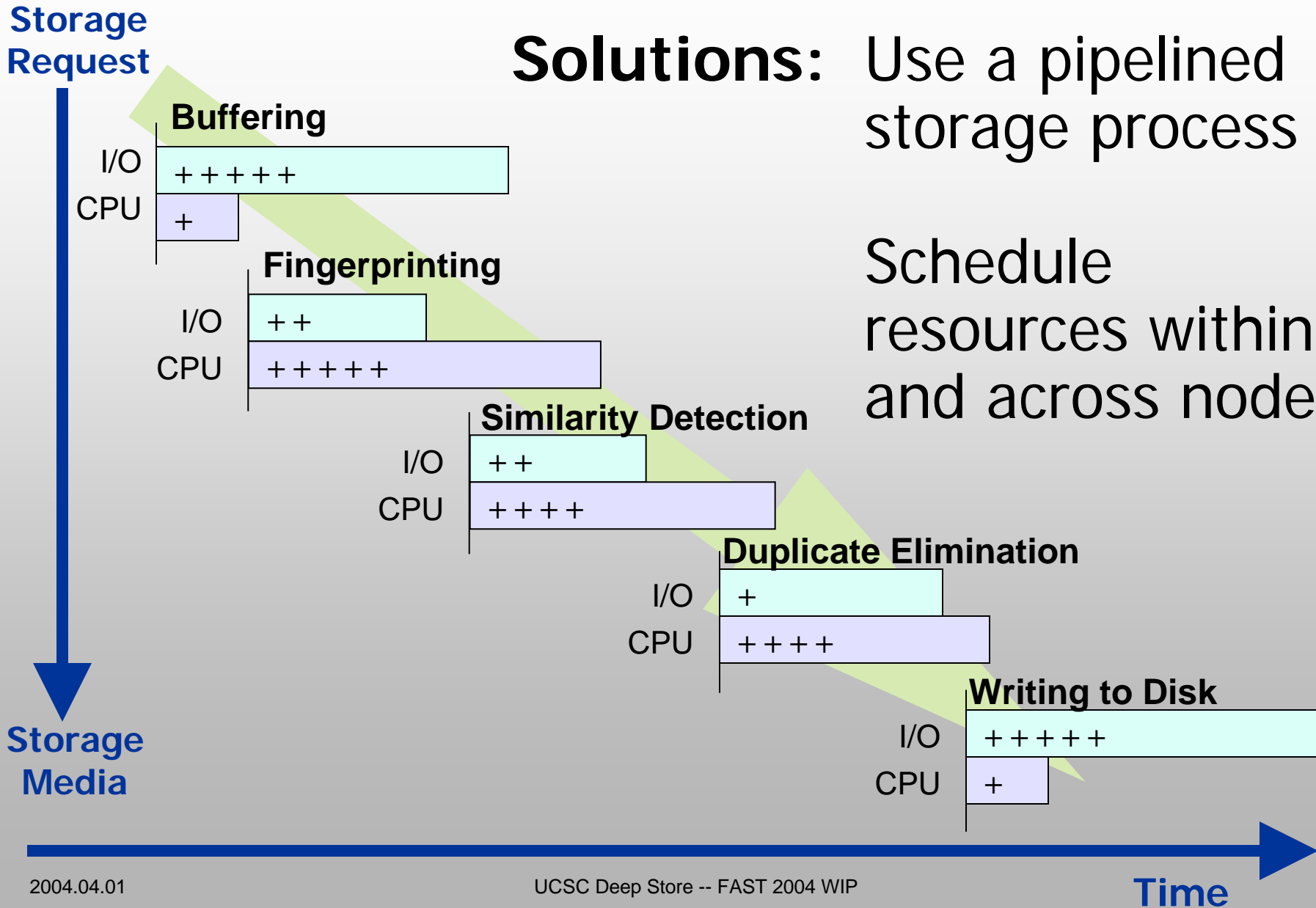**Problems:** The increasing size of content demands higher bandwidth

Users demand on-line behavior

Compression introduces additional costs to performance

# Storage Pipeline

**Solutions:** Use a pipelined storage process

Schedule resources within and across nodes



Storage Request

**Buffering**

| I/O | + + + + + |
| CPU | + |

**Fingerprinting**

| I/O | + + |
| CPU | + + + + + |

**Similarity Detection**

| I/O | + + |
| CPU | + + + + |

**Duplicate Elimination**

| I/O | + |
| CPU | + + + + |

**Writing to Disk**

| I/O | + + + + |
| CPU | + |

Storage Media
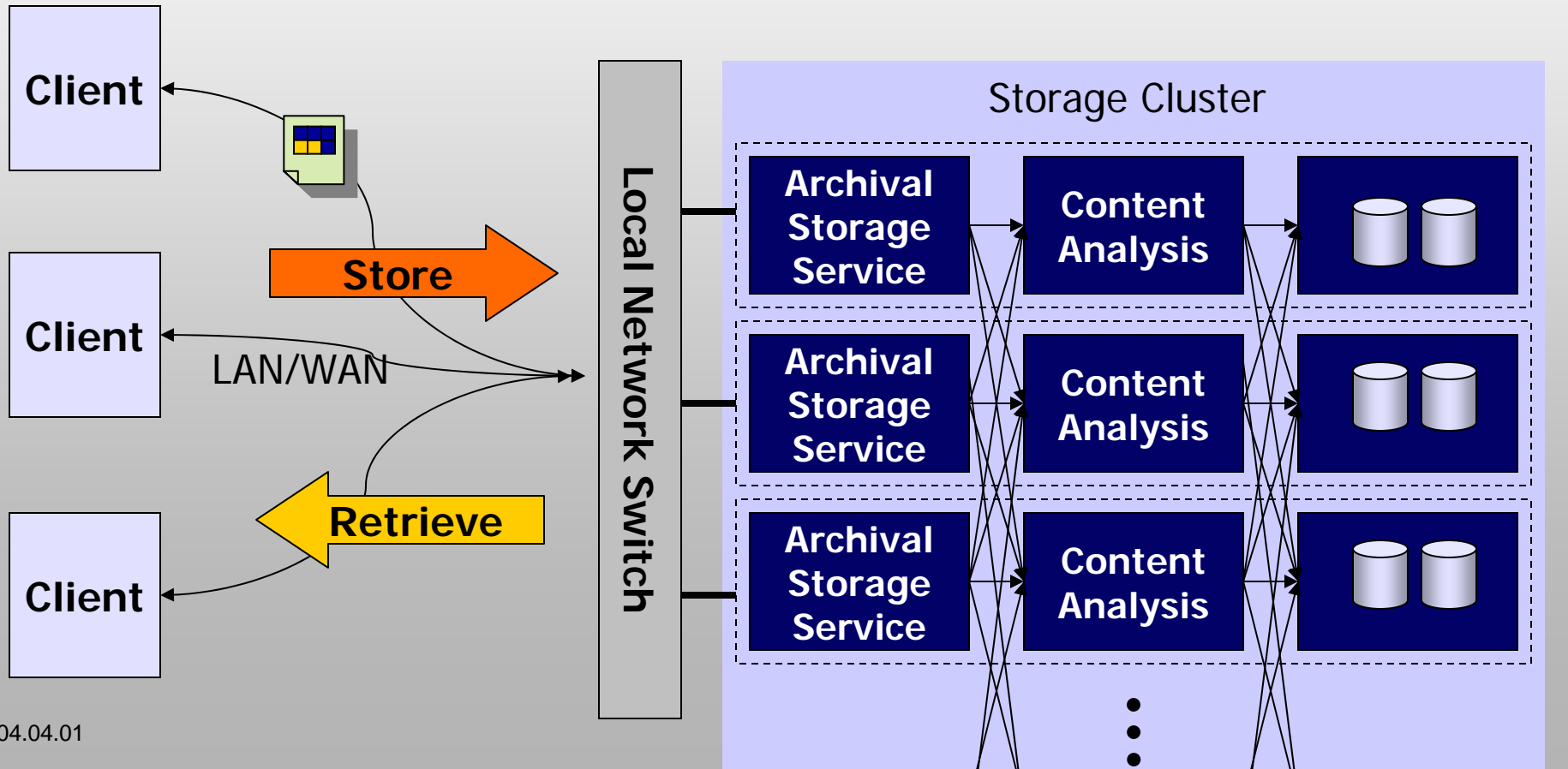
Time

# Problem 4: Managing Scale

**Problems:** Centralized terabyte to petabyte storage would create bottlenecks

Searching over all content is infeasible

# Distributed Archival Storage

**Solutions:** Use a distributed architecture

Reduce search space to metadata



2004.04.01

# Closing

▶ What is the Deep Store?

   A project developing an architecture and a working prototype to archive content on disk.

▶ Why is this different?

   Disk-based archival storage systems are not disk-based storage systems. These are different problems.

▶ How are we doing this?

   Design from the top down; build from the bottom up. We are developing an efficient, distributed node-based storage system.

▶ See our poster