

D

Data Longevity and Compatibility



Behrooz Parhami
Department of Electrical and Computer
Engineering, University of California, Santa
Barbara, CA, USA

Synonyms

[Data archives](#); [Data decay](#); [Media lifespan](#); [Storage failure](#)

Definition of Entry

Whether stored locally or in the cloud, data ultimately resides on physical media that are subject to electrical and physical deterioration and format obsolescence, making it necessary to augment the physical storage of data with a logical organization to protect data integrity and ensure its longevity.

Overview

Like many other attributes of computing and digital systems, the volume of data produced in the world is rising exponentially (Denning and Lewis 2016; Hilbert and Gomez 2011), with a growth rate that is even higher than those of

circuit density and processor performance, modeled by Moore's law (Brock and Moore 2006). A few exabytes of data generation per day in the early 2010s is slated to rise to many yottabytes in the 2020s (Cisco Systems 2017; Jacobson 2013). As data gains ever-greater value in the operation of social and business enterprises, data management, integrity, and preservation become major concerns. Any physical storage medium is subject to decay over time and has a finite lifespan. Some of these lifespans are relatively short compared with desirable data retention periods in practical settings. The format in which data is stored tends to become obsolete as well. It follows that an active strategy for ensuring the integrity and longevity of stored data is an important part of any data management plan.

Media for Long-Term Data Storage

The media used for storing data have undergone significant changes over time (Goda and Kitsuregawa 2012). Data resources are maintained in hierarchical arrangements, with hardware components in the hierarchy ranging from superfast (but expensive and, thus, low-capacity) caches near processing resources to vast (but relatively slow) data vaults with ultralow storage costs. The goal of the management scheme for such hierarchical arrangements is to make the data items of highest current value reside in faster storage, where they can be efficiently accessed,

with movements between the levels orchestrated so as to be responsive to changes in data values. Devices used for storing data for immediate access are known as hot storage media, in contrast to cold media that store data archives that are not in current use.

The order-of-magnitude correspondences shown in Table 1 are helpful in putting data volumes in perspective. In rough terms, the first four rows of the table (up to TBs) are currently within the domain of personal storage (although the numbers will no doubt expand in the coming years), whereas the last four are within the purview of large organizations, municipalities, or nation-states. At the rate data production is growing, it won't be long before we will have to decide whether to adopt proposed prefixes for 10^{27} and 10^{30} (Googology Wiki 2018)!

Data lifetimes range from seconds (for temporary results kept in scratchpad or working memory, as database transactions run to completion) to centuries or more (for historical archives). Each combination of data lifetime duration and data access requirement dictates the use of certain storage device types. Setting volatile semiconductor memories aside because of their nonpermanence, the most common storage technologies used in the recent past are listed in Table 2.

Data Decay and Device Lifespans

All storage media decay over time, although the degradation mechanism, and thus methods for dealing with it, is technology-dependent. Degradations that are small-scale and/or local can be compensated for through error-correcting codes. However, once a certain number of data errors have been corrected automatically, the storage medium must be discarded and a fresh copy of the data created in a different place, because continued deterioration will eventually exceed the codes' error-correcting capacity, leading to data loss.

Magnetic storage media, the most commonly used current technology for large-scale and long-term data repositories (hot or cold), degrade through weakening or loss of magnetization, a

process that is accelerated by external magnetic fields, heat, vibration, and a number of other environmental factors. Generally, the result of such deterioration is partial data erasure, rather than arbitrary changes to the stored information. Thus, when codes are used, they can be of the erasure variety (Plank 2013; Rizzo 1997), which imply lower redundancy compared with codes for correcting arbitrary errors. Archival magnetic media, such as reel tapes, must be stored under carefully regulated environmental conditions to maximize their lifespans.

Optical storage media can decay due to the breakdown of material onto which data is stored. Here too storing the media under proper environmental conditions can reduce data decay and increase the media's lifespan. When longer lifespans are desired, storage media of higher quality must be procured. In particular, M-discs and other specially developed archival media can prolong lifespans to many centuries (Jacobi 2018; Lunt et al. 2013). Such long lifespans can be validated only via accelerated testing (Svrcek 2009), a process whose results are not as trustworthy as those obtained through direct testing in actual field use.

Solid-state media use electrical charges to store data. Imperfect insulation can lead to charges leaking out and thus data being lost. Given that the device itself degrades much more slowly, refreshing data before total decay is one possible strategy for avoiding data loss. This process is in effect very similar to refreshing in DRAM chips (Bhati et al. 2016), except that the refreshing occurs at much slower rates and, thus, its performance hit and power waste are not serious issues.

Sources differ widely on the useful lifespans of data storage media. This is in part due to the fact that manufacturing quality associated with various suppliers and environmental conditions under which the media operate varies widely. Some typical figures for the most common media are included in Table 2. These figures should be used with caution, as there is no guarantee that even the lower ends of the cited ranges will hold in practice.

Data Longevity and Compatibility, Table 1 Developing an appreciation for data volumes

Abbr.	Name	Bytes	Example(s)
KB	Kilobyte	10 ³	One page of text or a small contact list
MB	Megabyte	10 ⁶	One photo or a short YouTube video
GB	Gigabyte	10 ⁹	A movie
TB	Terabyte	10 ¹²	Netflix's movies, 4 years worth of watching
PB	Petabyte	10 ¹⁵	Data held by an e-commerce site or bank
EB	Exabyte	10 ¹⁸	Google's data centers
ZB	Zettabyte	10 ²¹	WWW size or capacity of all hard drives
YB	Yottabyte	10 ²⁴	Worldwide daily data production by the 2020s

Data Longevity and Compatibility, Table 2 Attributes of storage devices in use over the past few decades

Type	Capacity	Cost	Speed	Life (years)	Archival	Additional information
Floppy disk	MB	Low	Low	5–7	No	Obsolete technology/format
CD/DVD	GB	Low	Medium	3–10	Yes	Are becoming less prominent
Blu-ray	GB+	Low	Medium	50	Yes	Still used for small-scale archives
M-disc	GB+	Low	Medium	1000	Yes	Lifespan unverified
Memory card	GB	High	High	5–10	No	Temporary personal storage
Cassette tape	GB	Low	Low	10–20	No	Obsolete technology/format
Cartridge	GB+	Low	Low	10–20	Yes	Used in mechanically accessed vaults
Hard disk	TB	Medium	Medium	3–5	No	Increasingly being replaced by SSDs
SSD (flash)	TB	High	High	5–10	No	Lifespan varies with write frequency
Reel tape	TB+	Low	Low	20–30	Yes	Lifespan with proper storage, low use
Disk array	PB	Medium	Medium	10–20	Yes	Lifespan improved by fault tolerance

Preventing, and Recovering from, Data Decay

Countermeasures for avoiding data decay are very similar in nature to those used to deal with data corruption, loss, or annihilation, discussed in this encyclopedia's entry "Data Replication and Encoding" and in general references on dependable and fault-tolerant computing (Parhami 2018). Low redundancy coding can counteract small-scale or local damages (Rao and Fujiwara 1989), while replication helps prevent widespread damage. Of course, replication comes with a high storage cost, which will be even more prohibitive in the age of big data. For this reason, data dispersion (Iyengar et al. 1998; Rabin 1989) and network coding (Dimakis et al. 2011) schemes, which in effect combine the advantages of low redundancy coding with the error correction strengths of replication, will be

the preferred methods for large collections of data.

Decades of experience with building and managing large database systems are being applied to the design of big data repositories that are nonvolatile and long life (Arulraj and Pavlo 2017). However, new challenges arising from enormous amounts of less structured and unstructured data remain to be overcome. Recovery methods (Arulraj et al. 2015) are both device- and application-dependent, so they need continual scrutiny and updating as storage devices evolve, and data is used in previously unimagined volumes and ways.

Future Directions

The importance of building large-scale, nonvolatile data archives in order to safeguard immense volumes of valuable data has been

recognized, to the extent of being discussed in the mass media (Coughlin 2014). Many projects are trying to achieve this goal (Yan et al. 2017). Some of the components needed for such system are already available and others are being invented.

On the storage technology side, disk arrays (Chen et al. 1994), which have played an important role in satisfying our data storage capacity and reliability needs for many years, must be adapted to the challenges presented by big data. We currently don't view semiconductor memories as suitable for data archives, but this assessment may change in light of technology advances (Qian et al. 2015).

On the data structuring, management, and maintenance side, data centers and their interconnection networks will play a key role in the development of robust permanent data archives (Chen et al. 2016). Improved understanding of storage device failure mechanisms and data decay (Petascale Data Storage Institute 2012; Schroeder and Gibson 2007) is another line that should be pursued.

Finally, development of new technologies for data storage, discussed in this encyclopedia under "Storage Technologies for Big Data" and "Emerging Hardware Technologies," will no doubt be accelerated as the storage and longevity needs to expand. Examples abound, but two promising avenues currently being pursued are storing vast amounts of data on DNA (Bornholt et al. 2016; Goldman et al. 2013) and building memories based on emerging nanotechnologies involving memristors and other new device types (Menon and Gupta 1999; Strukov et al. 2008).

Cross-References

- ▶ [Data Replication and Encoding](#)
- ▶ [Storage Hierarchies for Big Data](#)
- ▶ [Storage Technologies for Big Data](#)
- ▶ [Structures for Large Data Sets](#)

References

- Arulraj J, Pavlo A (2017) How to build a non-volatile memory database management system. In: Proceedings of the ACM international conference on management of data, Chicago, pp 1753–1758
- Arulraj J, Pavlo A, Dullloor SR (2015) Let's talk about storage & recovery methods for non-volatile memory database System. In: Proceedings of the ACM international conference on management of data, Melbourne, pp 707–722
- Bhati I, Chang M-T, Chishti Z, Lu S-L, Jacob B (2016) DRAM refresh mechanisms, penalties, and trade-offs. *IEEE Trans Comput* 65(1):108–121
- Bornholt J, Lopez R, Carmean DM, Ceze L, Seelig G, Strauss K (2016) A DNA-based archival storage system. *ACM SIGOPS Oper Syst Rev* 50(2):637–649
- Brock DC, Moore GE (eds) (2006) Understanding Moore's law: four decades of innovation. Chemical Heritage Foundation, Philadelphia
- Chen PM, Lee EK, Gibson GA, Katz RH, Patterson DA (1994) RAID: high-performance reliable secondary storage. *ACM Comput Surv* 26(2):145–185
- Chen T, Gao X, Chen G (2016) The features, hardware, and architectures of data center networks: a survey. *J Parallel Distrib Comput* 96:45–74
- Cisco Systems (2017) The zettabyte era: trends and analysis. White paper. <http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/vni-hyperconnectivity-wp.html>
- Coughlin T (2014) Keeping data for a long time, Forbes, available on-line at <http://www.forbes.com/sites/tomcoughlin/2014/06/29/keeping-data-for-a-long-time/#aac168815e26>
- Denning PJ, Lewis TG (2016) Exponential Laws of computing growth. *Commun ACM* 60(1):54–65
- Dimakis AG, Ramachandran K, Wu Y, Suh C (2011) A survey on network codes for distributed storage. *Proc IEEE* 99(3):476–489
- Goda K, Kitsuregawa M (2012) The history of storage systems. *Proc IEEE* 100:1433–1440
- Goldman N, Bertone P, Chen S, Dessimoz C, LeProust EM, Sipos B, Birney E (2013) Towards practical, high-capacity, low-maintenance information storage in synthesized DNA. *Nature* 494(7435):77–80
- Googology Wiki (2018) SI prefix. On-line document. http://googology.wikia.com/wiki/SI_prefix. Accessed 23 Feb 2018
- Hilbert M, Gomez P (2011) The World's technological capacity to store, communicate, and compute information. *Science* 332:60–65
- Iyengar A, Cahn R, Garay JA, Jutla C (1998) Design and implementation of a secure distributed data repository. IBM Thomas J. Watson Research Division, Yorktown Heights

- Jacobi J (2018) M-Disc optical media reviewed: your data, good for a thousand years. PCWorld. On-line document. <http://www.pcworld.com/article/2933478/storage/m-disc-optical-media-reviewed-your-data-good-for-a-thousand-years.html>
- Jacobson R (2013) 2.5 quintillion bytes of data created every day: how does CPG & retail manage it? IBM Industry Insights. <http://www.ibm.com/blogs/insights-on-business/consumer-products/2-5-quintillion-bytes-of-data-created-every-day-how-does-cpg-retail-manage-it/>
- Lunt BM, Linford MR, Davis RC, Jamieson S, Pearson A, Wang H (2013) Toward permanence in digital data storage. Proc Arch Conf 1:132–136
- Menon AK, Gupta BK (1999) Nanotechnology: a data storage perspective. Nanostruct Mater 11(8):965–986
- Parhami B (2018) Dependable computing: a multi-level approach, draft of book manuscript, available on-line at http://www.ece.ucsb.edu/~parhami/text_dep_comp.htm
- Petascale Data Storage Institute (2012) Analyzing failure data. Project Web site: <http://www.pdl.cmu.edu/PDSI/FailureData/index.html>
- Plank JS (2013) Erasure codes for storage stems: a brief primer. Usenix Mag 38(6):44–50
- Qian C, Huang L, Xie P, Xiao N, Wang Z (2015) Efficient data management on 3D stacked memory for big data applications. In: Proceedings of the 10th international design & test symposium, Dead Sea, pp 84–89
- Rabin M (1989) Efficient dispersal of information for security, load balancing, and fault tolerance. J ACM 36(2):335–348
- Rao TRN, Fujiwara E (1989) Error-control coding for computer systems. Prentice Hall, Upper Saddle River
- Rizzo L (1997) Effective erasure codes for reliable computer communication protocols. ACM Comput Commun Rev 27(2):24–36
- Schroeder B, Gibson GA (2007) Understanding disk failure rates: what does an MTTF of 1,000,000 hours mean to you? ACM Trans Storage 3(3):8, 31 pp
- Strukov DB, Snider GS, Stewart DR, Williams RS (2008) The missing memristor found. Nature 453(7191):80–83
- Svrcek I (2009) Accelerated life cycle comparison of Millenniata archival DVD. Final report for Naval Air Warfare Center Weapons Division, 75 pp
- Yan W et al (2017) ROS: a rack-based optical storage system with inline accessibility for long-term data preservation. In: Proceedings of the 12th European conference on computer systems, Belgrade, pp 161–174