



Abb. 1: Maschinensaal des LRZ, um 1980.

Langzeitarchivierung

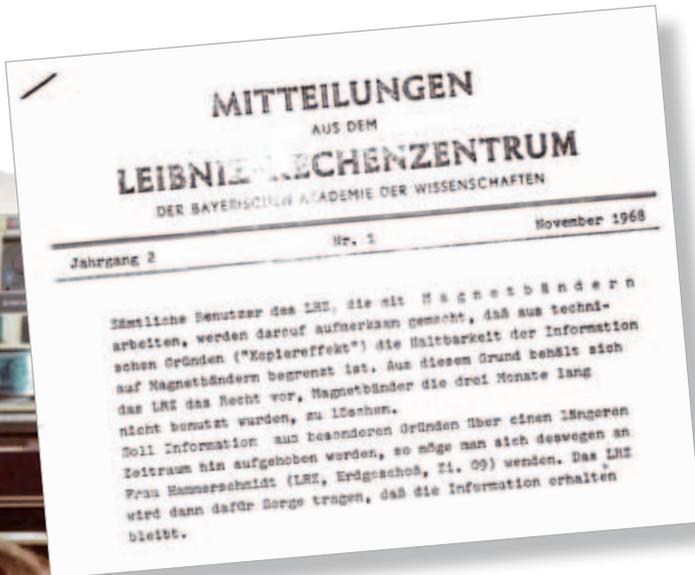
(K)Ein Grund zum Feiern

Die Geschichte des Leibniz-Rechenzentrums handelt nicht nur von Rechnern und Computern, sie belegt unter anderem auch die Geschichte von Daten und deren Speicherung: ein kurzer Rück- wie Ausblick auf die „unendliche“ Geschichte der langfristigen digitalen Archivierung von Daten im raschen technologischen Wandel der Zeit aus ungewohnter Perspektive.

VON WERNER BAUR

50 Jahre sind für ein Rechenzentrum eigentlich kein Grund zu feiern. 2, 4, 8, 16, 32 oder 64 sind runde Zahlen der Informationstechnologie. In solchen Einheiten werden auch die im Rechenzentrum gespeicherten Bytes gezählt. Apropos Bytes: Die ältesten Bytes im heutigen Archivsystem des LRZ sind genau 16 Jahre alt. Also doch ein Grund zum Feiern! Die Frage liegt nahe: Wie wurden Daten damals gespeichert, und wie ist es denn um die Archivierung in den nächsten 2, 4, 8 oder gar 16 Jahren bestellt? Eine Prognose für die nächsten 32 Jahre zu wagen, erscheint angesichts der rasanten Entwicklung der IT vermessen. Bleiben wir lieber bei 16, also im Jahr 2028, da ist der Autor hoffentlich in den Ruhestand gegangen und hat Zeit zurückzublicken.

Wir schreiben das Jahr 2028, Herr B. erinnert sich:



Robuster Datenträger der 1970er Jahre: die Lochkarte

Ende der 1970er Jahre hatte Herr B. in einer Informatik-Vorlesung an der TU München von der vor langer Zeit von Konrad Zuse entworfenen Rechenmaschine Z4 gehört, die die Nachkriegswirren im idyllischen Allgäu überstanden hatte und heute im Deutschen Museum in München steht. Die Z4 hatte einen mechanischen Speicher und konnte über einen Abtaster Lochstreifen lesen. Diese Maschine und ihre Zeit schienen Herrn B. unendlich weit weg von den modernen Computern am Leibniz-Rechenzentrum, der

Telefunken TR4, der TR440 oder gar den noch moderneren Cyber-Rechnern des Jahres 1979. Ihm wurde schnell klar, dass in der Informationstechnologie der Zeitbegriff „lang“ relativ „kurz“ zu verstehen ist, und da ihm diese Erkenntnis wichtig erschien, beschloss er, den Satz für die Nachwelt zu erhalten und zu archivieren. Zwar gab es schon seit über zehn Jahren Magnetbandspeicher an der Doppelprozessoranlage TR440 der Firma AEG-Telefunken und später dann an den Cyber-Rechnern am LRZ – als Erstsemesterstudent hatte Herr B. aber noch keinen Zugang zu den Magnetbandspeichern. Auch hatte er gehört, dass die Haltbarkeit von Magnetbändern seit Ende der 1960er Jahre zwar erheblich verbessert wurde, aber immer noch zu wünschen übrigließe (Abb. 2).

Das Mittel der Wahl war somit, den archivierungswürdigen Satz „Lang ist relativ kurz.“ binär zu kodieren und auf eine Lochkarte zu stanzen. Das hätte er am LRZ übrigens auch schon 15 Jahre vorher an der Telefunken TR4, die schon 1964 einen Lochkartenstanzer besaß, in der Richard-Wagner-Straße in München tun können (Abb. 3).

Die Lösung der 1980er Jahre: Magnetbandspeicher

Die Menge an Informationen, die eine Lochkarte speichern konnte, war mit 80 Byte recht begrenzt. Auch erwies sich die Lochkarte als keine dauerhafte Archivierungslösung. Zwar hätte das Medium selbst noch viele Jahre durchgehalten, aber die letzten Lochkartenleser wurden am LRZ 1985 abgeschafft. Herr B. musste sich um eine Ersatzlösung kümmern.

Abb. 2: Auszug aus den Benutzermittteilungen des Leibniz-Rechenzentrums, 1968.

Da Speicherplatz auf den Festplatten der Cyber-Großrechner immer noch sehr teuer war und andererseits die Zuverlässigkeit von Magnetbändern zwischenzeitlich deutlich gestiegen war, schien nun das 9-Spur-Magnetband mit einem Fassungsvermögen von 100 Megabyte die endgültige Lösung für das Archivierungsproblem des Herrn B. zu sein. So wurde das Zitat an den Cyber-Großrechnern des LRZ unter dem Betriebssystem NOS (Network Operating System), das die Magnetbandgeräte verwaltete, in einem speziellen Format auf ein Band geschrieben, wo es dann auch lange Jahre friedlich ruhte (Abb. 1 und 4).

Völlig neue Wege in den 1990ern: Bandroboter und Videokassetten

Anfang der 1990er Jahre wurde absehbar, dass auch die 9-Spur-Bänder keine endgültige Lösung waren. Die Großrechner sollten zugunsten von

Abb. 3: Lochkarte mit Stanzung „Lang ist relativ kurz“.



Abb. 4: 9-Spur-Bandgeräte und Magnetbandlager des LRZ.

Workstations und Vektorrechnern abgeschafft werden. Inzwischen war der Bedarf nach unabhängigen Archivsystemen erkannt worden. Mit der Installation des ersten Bandroboters am LRZ, betrieben unter Unitree, einer Software, die nach dem IEEE Mass Storage Reference Modell eigens für den Zweck der Archivierung von Massendaten entwickelt worden war, ging das LRZ völlig neue Wege. Zudem wurde ein neues Speichermedium eingesetzt: Videokassetten, wie sie auch in der Unterhaltungsindustrie genutzt wurden. Das Zitat des Herrn B. wurde vom 9-Spur-Band geholt und zusammen mit anderen Textfiles von den Datenträgern der sterbenden Großrechner auf eine Kassette des neuen Bandarchivs geschrieben.

1996 – Geburtsstunde des Archiv- und Backupsystems (ABS)

Auf den Videokassetten blieb den Bytes allerdings keine Zeit, um Staub anzusetzen. In der Theorie wohldurchdacht, erwies sich das System in der Praxis des täglichen Betriebs als äußerst fehleranfällig. Ein neuer Wechsel des Datenträgers und des gesamten Softwaresystems wurde notwendig. Diese so genannte Migration musste aber erstmalig nicht vom Nutzer des Archivs selbst durchgeführt werden. Den Job übernahmen die Betreiber des Systems im Rechenzentrum. Die Belegung im Archiv war inzwischen auf 1 Terabyte angewachsen, alle Daten der fehleranfälligen VHS-Videokassetten mussten mühselig ausgelesen und ins neue System geschrieben werden. Der Umzug der Daten zog sich über ein halbes Jahr hin. Dies war die Geburtsstunde des modernen Archiv- und Backupsystems (ABS), das auch 2012 noch am LRZ eingesetzt wurde. Anfang 1996 ging das neue System in Produktion. Als Datenträger wurden einspulige Kassetten in einem Bandarchiv der Firma IBM mit vier Bandlaufwerken eingesetzt. Eine völlig neue Semantik für den Begriff Bibliothek wurde geprägt. Zwar verstand der Programmierer am LRZ unter einer Bibliothek schon immer etwas anderes als etwa ein Bibliothekar an der Bayerischen Staatsbibliothek, nämlich eine Sammlung von Computer-Programmen und nicht eine Sammlung von Büchern. Aber dass eine Bibliothek auch eine Sammlung von Kassetten sein konnte, die von einem Roboter verwaltet wurde, war neu. Gesteuert wurde der Roboter von einem Softwarepaket, das damals noch ADSM



(Adstar Distributed Storage Manager) und später TSM (Tivoli Storage Manager) hieß und von IBM stammte. TSM führte genau Buch darüber, wo was in der Bandbibliothek lag. Außerdem übernahm TSM so lästige Arbeiten wie das Umkopieren der Daten im Archiv auf neue Datenträger. Dies war nicht etwa nötig, weil die Datenträger, also das Bandmaterial der Kassetten, schon nach wenigen Nutzungsjahren ermüdeten. Vielmehr erforderte der rasante Technologiefortschritt den Wechsel: Es gab wie einst bei den Lochkarten keine Lesegeräte mehr. So wanderte das Zitat von Herrn B. in den nächsten Jahren von den IBM 3590-Kassetten (Fassungsvermögen einer Kassette: 10 Gigabyte) auf STK Redwood Kassetten (50 Gigabyte) und von dort auf IBM LTO2-Kassetten (200 Gigabyte).

2004 – erste Archivdaten von der Bayerischen Staatsbibliothek

Neben dem regelmäßig notwendigen Umkopieren wurde ein anderes Problem immer präsenter. Es genügte nicht, dass eine Software namens TSM Buch darüber führte, auf welchem Band welche Datei lag. Auch Angaben über das Dateiformat, den Urheber, den Inhalt usw., also die so genannten Metainformationen, waren nötig, um die Information später wiederfinden zu können. Das Problem war nicht neu, die Archivare einer „richtigen“ Bibliothek kannten es schon lange. In jenen Jahren wurde die Zusammenarbeit des LRZ mit der Bayerischen Staatsbibliothek immer intensiver. TSM konnte die fehlenden Funktionalitäten auf der Anwenderseite nicht erbringen. Da-

für war es nicht gedacht. Neue Systeme wurden an TSM angebunden, die sich am in Bibliothekskreisen wohlbekannten OAIS-Referenzmodell (Open Archival Information System) orientierten. Nach einigen Jahren Ruhe auf LTO2-Kassetten war es 2010 wieder einmal Zeit für eine Verlagerung des Zitats auf neue Datenträger. Von den LTO2-Kassetten wanderte es auf die moderneren STK T10000-Kassetten (500 Gigabyte pro Kasette). Wie schon zuvor bekam auch diesmal der Autor davon nichts mit. Das erledigte TSM für ihn. Über 10.000 LTO2-Kassetten wurden nach Abschluss der Verlagerung am LRZ entsorgt. Aus Datenschutzgründen konnten die Kassetten nicht einfach weggeworfen werden. Sie wurden „zertifiziert“ verschrottet (Abb. 5).

attraktiv machte, trotz der hohen Kapazitäten von Festplatten und der Geschwindigkeit von Solid State Speicher. Immer noch wurden Archivdaten am LRZ deswegen auf Bändern gespeichert. Die großen Hersteller hatten gerade die neueste Generation der LTO-Reihe, LTO8, mit einem Fassungsvermögen von 12 Terabyte pro Band angekündigt. Die Archivdatenmenge der Staatsbibliothek von 2012, für deren Speicherung damals noch viele hundert Kassetten notwendig waren, hätte 2020 auf ein paar Dutzend LTO8-Kassetten Platz gehabt, aber natürlich hatte sich der Datenbestand in den letzten Jahren weiter vervielfacht.

Die jüngst eingeführte neue Version 7.0 von Rosetta, die seit einigen Jahren landesweit die einheitliche Schnittstelle zu den Archiven des LRZ bildete, etablierte sich nun bundesweit. Auch die ungeheuren Mengen an Forschungsdaten wurden nun nicht mehr direkt nach TSM geschrieben, sondern nutzten ein Rosetta-Derivat, das speziell für große Datenmengen geeignet war.

2028 – die Ära der Bänder geht zu Ende

Ein letztes Mal war das Zitat vor einigen Jahren auf eine neue Kassettengeneration umkopiert worden. Nun werden auch die letzten Bandgeräte am LRZ abgeschafft, die Daten werden auf die neuen Nanospeicher transferiert. Datenspeicher-Bibliotheken gibt es aber immer noch. Sie werden nun für die Nanospeichereinheiten genutzt. Die Bibliotheken füllen inzwischen den ersten Stock aller drei Datenwürfel des LRZ. Der dritte Würfel wurde erst vor einigen Jahren gebaut, da der Platz für die Speichermengen des europäischen 3-D-Datenarchivs nicht mehr ausreichte.

Für Herrn B. wird es langsam Zeit, sich darum zu kümmern, wo sein schöner Satz nach seinem Ableben bleibt. Aus sentimentalen Gründen hatte er sowieso die Lochkarte von einst aufbewahrt. Auch hatte er längst die Datei mit Metadaten versehen, ein SIP (submission information package) geschnürt und das Ganze in Rosetta 14.0 gepackt. Nun war das Zitat zu jeder Zeit und von überall aus der European Archive Cloud abrufbar. Außerdem gab es mehrere Kopien der Archivdatei an verschiedenen Standorten in Europa. Trotzdem – hatte er genug vorgesorgt für eine wirkliche Langzeitarchivierung?

Sicherheitshalber meißelt er das Zitat noch auf einen Stein in seinem Garten. Die Archivierung für die nächsten 1.000 Jahre sollte damit gesichert sein, auch wenn das dann immer noch keine unendliche Geschichte ist.

Abb. 5: Zertifizierte Entsorgung von LTO2-Kassetten.

2012 – ABS feiert 16. Geburtstag

Im Jahr 2012 hatte der Archivdatenbestand am LRZ beachtliche Ausmaße erreicht. Allein die Archivdaten der Bayerischen Staatsbibliothek umfassten 400 Terabyte, was übrigens mehreren Billionen Lochkarten entspricht. Absoluter Spitzenreiter waren aber die Massendaten der Supercomputer am LRZ mit 10.000 Terabyte. Anders als bei den Archivdaten der Staatsbibliothek, die über spezielle Frontends (ZEND, Digitool) ins Archiv gekommen waren, wurde für das Einbringen der Massendaten ins Archiv, den so genannten Ingest, TSM direkt genutzt.

2012 war auch das Jahr, in dem die Verbundzentrale des Bibliotheksverbunds Bayern die Archivierungssoftware Rosetta 3.0 einführte. Rosetta war „a new way of preserving cultural heritage and cumulative knowledge“, so die Herstellerfirma ExLibris.

2020 – Energiekosten halten das Band am Leben

Im IT-Geschäft war inzwischen Energie zum Kostenfaktor Nummer 1 geworden, eine Tatsache, die das Band als Speichermedium nach wie vor

DER AUTOR

Werner Baur leitet die Gruppe Datei- und Speichersysteme am Leibniz-Rechenzentrum der Bayerischen Akademie der Wissenschaften.