

10

Seismic Data Formats, Archival and Exchange

Bernard Dost, Jan Zednik, Jens Havskov, Raymond Willemann and
Peter Bormann

10.1 Introduction

Seismology entirely depends on international co-operation. Only the accumulation of large sets of compatible high quality data in standardized formats from many stations and networks around the globe and over long periods of time will yield sufficiently reliable long-term results in event localization, seismicity rate and hazard assessment, investigations into the structure and rheology of the Earth's interior and other priority tasks in seismological research and applications.

For almost a century, only parameter readings taken from seismograms were exchanged with other stations and regularly transferred to national or international data centers for further processing. Because of the uniqueness of traditional paper seismograms and lacking opportunities for producing high-quality copies at low cost, original analog waveform data, cumbersome to handle and prone to damage or even loss, were rarely exchanged. The procedures for carefully processing, handling, annotating and storing such records have been extensively described in the 1979 edition of the Manual of Seismological Observatory Practice (Willmore, 1979) in the chapter *Station operation*. They are not repeated here. Also the traditional way of reporting parameter readings from seismograms to international data centers such as the U.S. Geological Survey National Earthquake Information Center (NEIC), the International Seismological Centre (ISC) or the European Mediterranean Seismological Centre (EMSC) are outlined in the old Manual in detail in the section *Reporting output*. They have not changed essentially since then. On the other hand, respective working groups on parameter formats of the IASPEI and of its regional European Seismological Commission (ESC) have meanwhile debated for many years how to make these formats more homogeneous, consistent and flexible so as to better accommodate also other seismologically relevant parameter information.

Any data report, of course, must follow a format known to the recipient in order to be successfully parsed. Some of the goals for any format are:

- *concise* avoiding unnecessary expense in transmission and storage;
- *complete* providing all of the information required to use the data;
- *transparent* easily read by a person, perhaps without documentation; and
- *simple* straightforward to write and parse with computer programs.

10. Seismic Data Formats, Archival and Exchange

Traditional formats for reporting parameter data sacrificed simplicity, transparency and even sometimes completeness in favor of the other goals. With the falling cost of data storage and exchange, modern formats more often sacrifice conciseness in favor of transparency and simplicity.

In addition, modern formats are usually extensible and include “metadata”. An extensible format includes some way for new types of data to be introduced without either collecting all the new information into unformatted comment strings or making messages with the new data types unreadable by old parsers. “Metadata” are information about the data, such as how and by whom the data were prepared.

The Telegraphic Format (TF), as documented in the Manual of Seismic Observatory Practice (Willmore, 1979), is an extreme example of a traditional format for reporting and exchanging parameter data. Since telex was very expensive compared with modern communication costs, conciseness was the paramount goal even to the point of occasional ambiguity. The year of the data, for example, might be excluded if the recipient could probably infer it. The format was intended for use in an era when many stations were isolated and could report little more than their own phase readings, so event parameters such as hypocenter and magnitude were relegated to a secondary role. The TF incorporated further restrictions due to the special limitations of telex messages, such as no lower-case letters and sometimes no control over line breaks.

A seismic network with modern, calibrated instruments can provide far more information than telegraphic format allows, while low-cost e-mail has eliminated the restrictions and high costs of telex messages. Consequently, since at least 1990 most seismic parameter data have been stored and exchanged in modern formats that are more complete, simpler and usually more transparent than the Telegraphic Format. Until recently, however, there was no generally accepted standard modern format. A major step forward in this direction was made by the Group of Scientific Experts (GSE) organized by the United Nations Conference on Disarmament. It developed GSE/IMS formats (see 10.2.4) for exchanging parametric seismological data in tests of monitoring the Comprehensive Nuclear-Test-Ban Treaty (CTBT) (see 10.2.4) which became popular also with other user groups. Seismological research, however, has a broader scope than the International Monitoring System (IMS) for the CTBT. Therefore, a new *IASPEI Seismic Format* (ISF), compatible with the IMS format but with essential extensions, has been developed and adopted by the Commission on Seismological Observation and Interpretation of the International Association of Seismology and Physics of the Earth’s Interior at its meeting in Hanoi, August 2001. It is the conclusion of a 16-year process seeking consensus on a new format and fully exploits the much greater flexibility and potential of E-mail and Internet information exchange as compared to the older telegraphic reports (see 10.2.5).

Digital waveform data, however, are nowadays by far the largest volume of seismic data stored and exchanged world-wide. The number of formats in existence and their complexity far exceeds the variability for parameter data. With the wide availability of continuous digital waveform data and unique communication technologies for world-wide transfer of such complete original data, their reliable exchange and archival has gained tremendous importance. Several standards for exchange and archival have been proposed, yet a much larger number of formats are in daily use. The purpose of the section on digital waveform data is to describe the international standards and to summarize the most often used formats. In addition, there will be a description of some of the more common conversion programs.

Beforehand, however, a short description of the most common parameter formats is given below.

10.2 Parameter formats

Parameter formats deal with all earthquake parameters like hypocenters, magnitudes, phase arrivals etc. Until recently, there were no real standards, except the Telegraphic Format (TF) used for many years to report phase arrival data to international agencies (Willmore, 1979; Chapter “Reporting output”). The format is not used for processing. There have been attempts to modernize TF for many years through the IASPEI Commission of Practice (now the Commission on Seismological Observation and Interpretation) and as mentioned in the introduction, the IASPEI Seismic Format (ISF) was approved as a standard in 2001. In practice, many different formats are used and the most dominant ones have come from popular processing systems. In the following, some of the most well known formats will be briefly described. For complete description of the formats, the reader is referred to original Manuals or publications.

10.2.1 HYPO71

The very popular location program HYPO71 (Lee and Lahr, 1975) has been around for many years and has been the most used program for local earthquakes. The format was therefore limited to work with only a few of the important parameters. Tab. 10.1 gives an example.

Tab. 10.1 Example of an input file in HYPO71 format. Each line contains, from left to right: Station code (max 4 characters), E (emergent) or I (impulsive) for onset clarity, polarity (C – compression; D – dilatation), year, month, day, and time (hours, minutes, seconds, hundredth of seconds) for P-phase onset, second for S-phase onset (seconds and hundredth of seconds only), and, in the last column, record duration. The blank space between ES and duration has been used for different purposes like amplitude. The last line is a separator line between events and contains control information.

FOO	EPC	96	6	6	64848.47	62.67ES	136
MOL	EPC	96	6	6	64849.97	65.87ES	144
HYA	EP	96	6	6	64856.78	78.07ES	135
ASK	EP	96	6	6	649 2.94	34.72ES	183
BER	EPC	96	6	6	649 7.56	36.61ES	
EGD	EPD	96	6	6	649 5.76	40.53ES	
					10 5.0		

The format is rather limited since only P- or S-phase names can be used and the S phase is referenced to the same hour-minute as the P phase; also, the format can not be used with teleseismic data. However, it is probably one of the most popular formats ever for local earthquakes. The HYPO71 program has seen many modifications and the format exists in many forms with small changes.

10.2.2 HYPOINVERSE

Following the popularity of HYPO71, several other popular location programs followed like Hypoinverse (Klein, 1978) and Hypoellipse (Lahr, 1989). Tab. 10.2 gives an example of the input format for Hypoinverse.

Tab. 10.2 Example of the Hypoinverse input format. Note that year, month, day, hour, min is only given in the header and only one phase is given per line.

```
96 6 60648
FOO EPC 48.5 136
FOO ES 62.7
MOL EPC 50.0 144
MOL EPC 50.9
MOL ES 65.9
```

10.2.3 Nordic format

In the 1980's, there was one of the first attempts to create a more complete format for data exchange and processing. The initiative came from the need to exchange and store data in Nordic countries and the so-called Nordic format was agreed upon among the 5 Nordic countries. The format later became the standard format used in the SEISAN data base and processing system and is now widely used. The format tried to address some of the shortcomings in HYPO71 format by being able to store nearly all parameters used, having space for extensions and useful for both input and output. An example is given in Tab. 10.3.

Tab. 10.3 Example of Nordic format. The data is the same as seen in Tabs. 10.1 and 10.2. The format starts with a series of header lines with type of line indicated in the last column (80) and the phase lines are following the header lines with no line type indicator. There can be any number of header lines including comment lines. The first line gives among other things, origin time, location and magnitudes, the second line is the error estimate, the third line is the name of the corresponding waveform file and the fourth line is the explanation line for the phases (type 7). The abbreviations are: STAT: Station code, SP: component, I: I or E, PHAS. Phase, W: Weight, D: polarity, HRMM SECON: time, CODA: Duration, AMPLIT: Amplitude, PERI: Period, AZIMU: Azimuth at station, VELO: Apparent velocity, SNR: Signal-to-noise ratio, AR: Azimuth residual of location, TRES: Travel-time residual, W: Weight in location, DIS: Epicentral distance in km and CAZ: Azimuth from event to station.

```
1996 6 6 0648 30.4 L 62.635 5.047 15.0 TES 13 1.4 3.0CTES 2.9LTES 3.0LNAO1
GAP=267 5.92 18.8 43.0 31.8 -0.5630E+03 0.8720E+03 -0.3916E+03E
1996-06-06-0647-46S.TEST__011 6
STAT SP IPHASW D HRMM SECON CODA AMPLIT PERI AZIMU VELO SNR AR TRES W DIS CAZ7
FOO SZ EP C 648 48.47 136 -0.110 116 180
FOO SZ ESG 649 2.67 0.710 116 180
FOO SZ E 649 2.89 426.4 0.3 116 180
MOL SZ EP C 648 49.97 144 -0.310 129 92
MOL SZ EPG C 648 50.90 0.410 129 92
MOL AZ E 649 5.86 129 92
MOL SZ ESG 649 5.87 0.410 129 92
MOL SZ E 649 6.98 328.6 0.6 129 92
HYA SZ EP 648 56.78 135 0.810 174 159
HYA SZ IP D 648 56.78 0.810 174 159
HYA SZ EPG D 648 57.56 0.110 174 159
```

```

HYA SZ ESG          649 18.07                                0.610 174 159
NRA0 SZ Pn          0649 24.03                                309.6  8.5 139  5 -0.410 403 119
NRA0 SZ Pg          0649 32.60                                305.6  7.285.2  1  0.410 403 119
NRA0 SZ Lg          0650 22.05                                302.0  4.016.0 -1 -0.410 403 119

```

10.2.4 The GSE/IMS formats

The GSE format (versions GSE1.0 and GSE2.0) was originally developed by the Group of Scientific Experts (GSE) of the Conference on Disarmament in Geneva and was used for the global technical test GSETT-3 organized by the GSE. With the establishment of the International Monitoring System (IMS) for the Comprehensive Nuclear-Test-Ban Treaty (CTBT) monitoring a significantly revised version of this format, termed GSE 2.1, was renamed to IMS1.0. This format has been widely used by many institutions around the globe, particularly in AutoDRM data exchanges (<http://seismo.ethz.ch/autodrm>) and for data transmission to international data centers, however less as a processing format than HYPO71 or the Nordic format. IMS1.0 is similar in structure to the Nordic format but more complete in some respects and lacking features in other. A major difference is that the line length can be more than 80 characters long, which is not the case for any of the previously described formats. After SEISAN, IMS1.0 is the first major format for which completeness or readability has been recognized as a more important design goal than conciseness.

The official custodian of the IMS format is the Comprehensive Nuclear-Test-Ban Treaty Organisation (CTBTO). As of December 2002, 166 States signed the CTBT and are participating in the development of the IMS system. The WEB page of CTBTO is <http://www.ctbto.org>. The IMS1.0 data format description can be obtained through National Data Centres (NDC) for CTBT which have been established in many countries on all continents. It is also available from the web site of the former Prototype International Data Centre (PIDC) under the heading "3.4.1 Rev3 Formats and Protocols for Messages" via <http://www.cmr.gov/pidc/librarybox/idcdocs/idcdocs.html>. It can be expected that in future CTBTO will post on its WEB page updates of its data formats, including the IMS format.

Tab. 10.4 Example of the IMS1.0 parameter format which contains the same data as given in Tabs. 10.1 to 10.3. The first lines are message information etc. The remaining lines are more or less self-explanatory. Note that more information, with a higher accuracy, can be given for each phase (like magnitude) than in the Nordic format. On the other hand, information like component and event duration is missing. These are added in the new ISF format.

```

BEGIN GSE2.0
MSG_TYPE DATA
MSG_ID 1900/10/19_1711 ISR_NDC
DATA_TYPE ORIGIN GSE2.0
EVENT 00000001

```

Date	Time	Latitude	Longitude	Depth	Ndef	Nsta	Gap	Mag1	N	Mag2	N	
1996/06/06	06:48:30.4	62.6350	5.0470	15.0	25	13	267			ML 2.9	8	
1.40	+- 5.92	0.0	0.0	0	+- 31.8	1.04	4.84			+-0.3		
Sta	Dist	EvAz	Phase	Date	Time	TRes	Azim	AzRes	Slow	SRes	Def	SNR
FOO	1.04	180.0	mc P	1996/06/06	06:48:48.5	-0.1					T	
FOO	1.04	180.0	m SG	1996/06/06	06:49:02.7	0.7					T	
FOO	1.04	180.0	m	1996/06/06	06:49:02.9							
426.4	0.30	ML 3.2		00000003	(from previous line)							
MOL	1.16	92.0	mc P	1996/06/06	06:48:50.0	-0.3					T	
MOL	1.16	92.0	mc PG	1996/06/06	06:48:50.9	0.4					T	

10. Seismic Data Formats, Archival and Exchange

```
MOL      1.16  92.0 m           1996/06/06 06:49:05.9
MOL      1.16  92.0 m   SG     1996/06/06 06:49:05.9   0.4                T
MOL      1.16  92.0 m           1996/06/06 06:49:07.0
NRA0     3.62 119.0 m   Pn     1996/06/06 06:49:24.0  -0.4 309.6   5.0  8.5   TAS  13.9
(from previous line)
NRA0     3.62 119.0 m   Pg     1996/06/06 06:49:32.6   0.4 305.6   1.0  7.2   TAS  85.2
NRA0     3.62 119.0 m   Lg     1996/06/06 06:50:22.0  -0.4 302.0  -1.0  4.0   TAS  16.0
(from previous line)
STOP
```

10.2.5 The IASPEI Seismic Format (ISF)

The need for an agreed-upon parameter format for comprehensive seismological data exchange has led to the IASPEI Seismic Format (ISF), adopted as standard in August 2001. ISF conforms to the IMS.1.0 standard but has essential extensions for reporting additional types of data. This allows the contributor to include complementary data considered to be important for seismological research and applications by the IASPEI Commission on Seismological Observation and Interpretation. The format looks almost like the IMS1.0 example in Tab. 10.4 above, except for the extensions. The ISF has been comprehensively tested at the ISC and NEIC and incompatibilities have been eliminated. The definite detailed description of the ISF is available from the ISC home page and kept up-to-date there (see <http://www.isc.ac.uk/Documents/isf.pdf>). Therefore, it is not reproduced in this Manual.

Consensus on the ISF was reached partly by including many optional items, so the format is not as simple as some alternatives. Despite this, the completeness, transparency, extensibility and metadata of ISF are expected to make it very widely used. Wide use of ISF will bring back the advantages of a generally accepted standard so that it becomes easier to exchange data, re-use data collected for past projects, and employ programs developed elsewhere.

In Volume 2, IS 10.1 and IS 10.2, examples are given of how event parameter data and unassociated parameter readings by seismic stations are reported according to the IMS format with ISF extensions.

10.3 Digital waveform data

Many different formats for digital data are used today in seismology. For a summary and the abbreviations used, see the following sections. Most formats can be grouped into one of the following five classes:

- 1) local formats in use at individual stations, networks or used by a particular seismic recorder (e.g., ESSTF, PDR-2, BDSN, GDSN);
- 2) formats used in standard analysis software (e.g., SEISAN, SAC, AH, BDSN);
- 3) formats designed for data exchange and archiving (SEED, GSE);
- 4) formats designed for database systems (CSS, SUDS);
- 5) formats for real time data transmission (IDC/IMS, Earthworm).

Use of the term "designed" in describing Class 3 and 4 formats is intentional. It is usually only at this level that very much thought has been given to the subtleties of format structure which result in efficiency, flexibility and extensibility.

The four classes (1-4) show a hierarchical structure. Class 4 forms a superset of the others, meaning that classes 1-3 can be deduced from it. The same argument applies to class 3 with respect to classes 1 and 2. Nearly all format conversions performed at seismological data centers are done to move upwards in the hierarchy for the purpose of data archiving and exchange with other data centers. Software tools are widely available to convert from one format to another and particularly upwards in the hierarchy.

This hierarchy also explains why there are so many formats. The design of class 1 formats depends on the manufacturer of the data acquisition system. In the early days of digital seismometry, display and analysis software was often proprietary and marketed specifically for a certain manufacturer's equipment and data format. There was no real need for manufacturers to adhere to a standard recording format, until users began to realize the advantages of exchanging data with other seismologists and discovered that this was quite difficult unless the other party was using the same hardware and/or software.

Station operators, who were not satisfied with the proprietary analysis software supplied with the procured data acquisition systems, started to convert data from Class 1 formats into the Class 2 formats which were used by more powerful and widely available analysis packages such as SAC. These programs usually provide subroutines that make conversion from local formats fairly easy. New analysis packages (e.g., SeisGram) which are developed around a Class 1 format (BDSN in this case) implicitly offer their format preference as a candidate for a new standard in Class 2, but it hardly matters as long as the necessary software tools are available to convert to and from the data exchange formats.

The GDSN (Global Digital Seismic Network) format began as a Class 1 format, but because it was used by an important global seismograph network (DWWSSN, SRO), it became accepted as a de facto standard for data exchange (Class 3). The beginning of widespread international data exchange within the FDSN (Federation of Digital Seismic networks) and GSE (Group of Scientific Experts) groups in the late 1980s revealed the GDSN format's weaknesses in this role and put in motion the process of defining more capable exchange formats.

The volume of commonly available digital seismic data continues to increase dramatically. It increased from 600 MB annually in 1980 to 300 GB in 1992 and today we are talking about many terabytes. Database systems, which are specially designed to handle these large datasets, have therefore begun to appear as a superset of the standard data exchange formats. The SUDS system is an example of this type of format.

In the 1990s, several activities (e.g., the GSETT-3 experiment and the U.S. National Seismograph Network (USNSN) have emerged which feature real-time exchange of seismological data, and interest has focused on formats which are suitable for such applications. In the late 1990s, this idea was carried farther by systems such as Earthworm, which implement format-independent protocols. Earthworm also is designed to exchange data across a peer network of multiple, independent nodes, as well as in a traditional network of dependent nodes with a centralized collection and distribution center.

Following is a brief description of some of the classes of formats as defined above.

10.3.1 Data archival

Data archival requires the storage of complete information on station, channel(s) and the structure of the data. Most existing formats are designed to provide part of the information. Most archival formats presently in use do include information on station and channel, but are not always complete in the description of the data. What we envisage is demonstrated through several features in the Standard for the Exchange of Earthquake Data (SEED) format:

- Data Description Language (DDL)
- reference to byte order;
- response information

The DDL is defined to enable the data itself to be stored in any data format (integer, binary, compressed). The language consists of a number of keys defining, for example, the applied compression scheme, number of bytes per sample, mantissa and gain length in bits and the use of the sign convention. The reader interprets the DDL and knows exactly how to deal with the data. The advantage of the DDL is that the original data structure can be maintained and is known. A disadvantage is that readers will have to interpret the DDL and have less performance in reading. However, the decoding information is available directly with the data and this is extremely important, since data are collected on platforms having different byte orders. In SEED the byte order of the original data is defined in the header, so the reader will be able to decide whether the data should be swapped.

In most archival formats, response information can be supplied in terms of poles and zeroes. Fewer efforts are undertaken to give the FIR filter coefficients in the header, although they are accounted for in the definition of SEED and GSE2.X. A problem occurs when a description of the instrument response is given only in measured amplitude and phase data as a function of frequency, as is the case in the GSE1.0 format. Also, the GSE2.X does not specify what is a minimum requirement. The main purpose of the response information is to correct for instrument response and thus the user will have to find the best fitting poles and zeroes to the given response. Although tools are available to calculate poles and zeroes from frequency, amplitude and phase data (e.g., in Preproc), results from the multiple inversion of the discrete frequency, amplitude and phase data will be different from the original data.

The deployment of large mobile arrays consisting of heterogeneous instrumentation is an important research tool. Data archival of these data is important. Although there is a tendency to store the data in a common format, the responses of sensors and data acquisition systems are often poorly known. It is recommended to pay attention to this issue before the experiment starts!

Finally, an issue in data archival is the responsibility of the data quality and the mechanism of reporting data errors. The network/station operator is responsible for the quality of the original data. However, the data may be subjected to format conversion at a remote data center. This last stage could introduce errors and it is the originator of the data, which must be responsible for data quality and should agree on the final conversion, if such a conversion is done externally.

10.3.2 Data exchange formats

The data exchange formats are closely related to the way data is exchanged. Therefore, these formats are described separately. Essentially, any format can be used for exchange, however the idea of an exchange format is to make it easy to send electronically, have a minimum standard of content and be readable on all computer platforms.

At present, there are many different techniques in use to exchange data, either between data users and data centers or between data centers. An overview of existing techniques is given below.

	Technique	Advantage	Disadvantage
Indirect on-line	autoDRM, NetDC	email based (no connection time)	small volume or download through ftp
Direct on-line	ftp, WWW, DRM (Spyder/Wilber/FARM)	direct access, enables easy data selection	slow for large data volumes
Off-line	CD-ROM (DVD)	direct access	no real-time data

Indirect on-line data exchange is arranged through (automated) *Data Request Managers* (DRMs) where the request mechanism is based on email traffic. There is work towards standardization on AutoDRM (<http://seismo.ethz.ch/autodrm>) to prevent a situation where users will have to learn a multitude of data request mechanisms with each having its own specific request format. One step further is the implementation of a communication protocol for exchange between data centers in such a way that a user only has to send one request to a nearby data center node. His/her request is then automatically routed through the data centers that may contribute to the requested data set. Such a protocol is under development and is known as the NetDC initiative (Casey and Ahern, 1996).

One basic problem in using email as the transport mechanism is the restricted data volume that can be exchanged. Also, the format sometimes will have to be ASCII. The format issue is taken care of in the GSE format, although in the description of the AutoDRM protocol it is mentioned that also a format like SEED can be used. The only difference is that the user is requested to get the data through anonymous ftp (pull) or the data is pushed into an anonymous ftp area defined by the user. The AutoDRM system at the Orfeus Data Centre (ODC) supports the SEED format in data exchange.

Direct on-line access to data is arranged at the ODC, for example, mainly through a website (<http://orfeus.knmi.nl>). A distinction is made between near real-time data collection (Spyder) and complete data volumes (ODC-volumes, FARM). Spyder data are available within a few hours after a major event, while ODC volumes lag behind real-time. At this moment there is a delay of approximately 3-4 years.

Internet speed is presently still limiting the usefulness of this direct on-line data exchange, especially since the volumes that are to be transferred may be large. One major advantage of direct on-line availability of the data is the capability to make a selection out of the vast amount of digital data. Procedures are presently under development to increase the power of these selection tools.

10. Seismic Data Formats, Archival and Exchange

Off-line data access provides complete, quality controlled data that are locally available at each institute in the form of CD-ROMs. The completeness and quality control takes time and CD-ROMs have a limited data volume. Digital Versatile Disks (DVDs) will probably replace CDs in the near future.

10.3.3 Formats for data base systems

Formats for data base systems are specially designed and no details will be given here. Examples of such formats are CSS and the derived “IDC Database Schema” (see IS 10.3 and <http://www.cmr.gov/pidc/librarybox/idcdocs/idcdocs.html>) and SUDS.

10.3.4 Continuous data protocols and formats

With better communication systems, real time transmission of digital data becomes more common. There is no internationally agreed upon format for this and equipment manufacturers use their own formats. The most widely used standard format is at present the CD-1.0 protocol used by the International Data Centre (IDC) for the International Monitoring System (IMS) as described under 10.2.4. Complete documentation can be found on the secure website <https://www2.ctbto.org> (authorized users only) and openly on <http://www.cmr.gov/pidc/librarybox/idcdocs/idcdocs.html>.

Up to 100 channels from a station or array of stations can be transmitted in near-real time using a single connection. Digital data are provided in compressed or uncompressed format and with or without authentication signatures. The protocol uses units of information called frames to establish or alter a connection and to exchange data between the sender and the receiver. Only one frame is being transmitted or received at any instance. A time-out is used in case of lost connection.

Establishing connections. The sender initiates the connection with the receiver to a pre-designated IP address and port by sending a Connection Request Frame. The receiver validates the authenticity of the sender and provides a new port and Internet Protocol (IP) address in a Port Assignment Frame. The sender drops the original connection and connects to the assigned IP address and port that is subsequently used for all data transfer.

Transmitting data. After the connection is established, the sender sends a Data Format Frame, which describes the format of the subsequent Data Frames. The sender can then send Data Frames data. The Data Format Frame provides information about itself and about Data Frames that will follow. The Data Frame contains the raw time series data. Each Data Frame has a single Data Frame Header and multiple channel sub-frames.

Altering connections. Either the sender or the receiver can alter the connection through the exchange of Alert Frames. The receiver sends the Alert Frame to notify the sender to use a different port. The sender uses Alert Frames to notify the receiver that the communication will cease or that a new data format is about to be used.

Terminating connections. Typically, an established connection remains active and in use until the sender or receiver terminates it for maintenance or reconfiguration. The connection can be intentionally terminated by sending an Alert Frame. Unintentional termination due to a slow or failed communications system is detected after the time-out period.

The CD-1.0 protocol is being replaced by the CD-1.1 protocol for transmission of IMS data; a description can be found on <https://www2.ctbto.org> and <http://www.cmr.gov/pidc/librarybox/idcdocs/idcdocs.html>.

Another real time data protocol is Earthworm, which is being used in North America. Documentation for this protocol can be found on the USGS website <http://gldbrick.cr.usgs.gov>.

10.4 Some commonly encountered digital data formats

Following is an alphabetical list of formats in use. For each format some description is given. The list of formats, of course, is not be complete, particularly for formats in little use, however, the most important formats in use today (2000) are included. In a later section, a list of popular analysis software systems is mentioned as well as a brief description of some conversion programs.

Only those formats are listed which can be converted by at least one of these analysis software systems. It is of particular importance to know on which computer platform the binary file has been written since only a few analysis programs work on more than one platform. Therefore, the data file should usually be written on the same platform as the one on which the analysis program is run. Accordingly, we will mention below, for each format, the respective computer platform.

AH

Class: 2 Platform: Unix

The Ad Hoc (AH) format is used in the AH waveform analysis software package developed at Lamont Doherty Geological Observatory, N.Y., USA. This package also supports a number of conversion tools.

CSS

Class: 2,4 Platform: Unix

The Center for Seismic Studies (CSS) Database Management System (DBMS) was designed to facilitate storage and retrieval of seismic data for seismic monitoring of test ban treaties [CSS]. The seismic data separate into two categories: waveform data and parametric data.

For the parametric data, the design utilizes a commercial relational database management system. Information is stored in relations that resemble flat, two-dimensional tables as in the ISF format (see annexed IS 10.1). The description of waveform data is physically separated from the waveform data itself. The index to the waveform archive is maintained within the relational database. Data are stored in plain files, called non-DBMS files. Each non-DBMS file is indexed by a relation that contains information describing the data and the physical location of the data in the file system. Each waveform segment contains digital samples from only one station and one channel. The time of the first sample, the number of samples and the sample rate of the segment are noted in an index record. The index also defines in which file

10. Seismic Data Formats, Archival and Exchange

and where in the file the segment begins, and it identifies the station and channel names. A calibration value at a specified frequency is noted. The index records are maintained in the **wfdisc** relation. Each **wfdisc** record describes a specific waveform segment and contains an id number to designate detailed information on the station and instrumentation of the trace.

GeoSig

Class: 1 Platform: PC

Binary format used by GeoSig recorders. The format consists of a header and multiplexed data.

Güralp format

Class: 1 Platform: PC

Format used by Güralp recorders.

ESSTF binary

Class: 1 Platform: All

The European Standard Seismic Tape Format (ESSTF) grew out of a major corporate effort by Lennartz Electronic GmbH [LEN]. ESSTF has been used as the framework for the file system in the SAS-58000 data acquisition system. ESSTF combines header information in ASCII format with seismic data in binary format.

The event header block is a single block preceding the data blocks, containing information on event start time. Each data block contains a 48-character header block (channel number, time, etc.) in ASCII. All channels are stored in a multiplexed form in one file. Data are organized in frames, each containing 500 data points. The most efficient access to the binary data is by unformatted, buffered reading with the capability of decoding the ASCII data directly out of a memory buffer.

GSE

Class 3 Platform: All

The format proposed by the Group of Scientific Experts (GSE format) has been extensively used with the GSETT projects on disarmament. The GSE2.1, now renamed IMS1.0, is the most recent version. The manual can be downloaded from (http://orfeus.knmi.nl/manuals/provisional_GSE2.1.ps) or the web pages of the Center for Monitoring Research in Arlington (http://www.cmr.gov/web-gsett3/CRP-243/www/FmtProt/FmtProt_5.html#HEADING113; <http://www.cmr.gov/pidc/librarybox/idcdocs/idcdocs.html>).

A GSE2.1 waveform data file consists of a waveform identification line (WID2) followed by the station line (STA2), the waveform information itself (DAT2), and a checksum of the data (CHK2) for each DAT2 section (Provisional GSE 2.1 Message Formats & Protocols, 1997). The default line length is 132 bytes. No line may be longer than 1024 bytes. The response data type allows the complete response to be given as a series of response groups that can be cascaded. Response description is made up of the CAL2 identification line plus one or more of the PAZ2, FAP2, GEN2, DIG2 and FIR2 response sections in any order.

Waveform identification line WID2 gives the date and time of the first data sample; the station, channel and auxiliary codes; the sub-format of the data, the number of samples and sample rate; the calibration of the instrument represented as the number of nanometers per digital count at the calibration period; the type of the instrument, and the horizontal and vertical orientation.

Line STA2 contains the network identifier, latitude and longitude of the station, reference coordinate system, elevation and emplacement depth.

Data section after DAT2 may be in any of six different sub-formats recognized in the GSE2.1 waveform format: INT, CM6, CM8, AUT, AU6, and AU8. INT is a simple ASCII sub-format, "CM" sub-formats are for compressed data and "AU" sub-formats are for authentication data. All represent the numbers as integers and therefore can be sent by email.

A checksum CHK2 must be provided in the GSE2.1 format. The checksum is computed from integer data values prior to converting them to any of the sub-formats.

IRIS dial-up expanded ASCII

Class: 1 Platform: All

The IRIS dial-up data retrieval system can be used to search for, display, and write data from IRIS GSN stations which are equipped with dial-up capabilities. Digital waveforms can be written in ASCII using the various on-line commands, e.g., "V" variable- and "F" fixed-record-length, expanded ASCII. These files contain two types of records: header records (one per file) and data records. The header record contains station and instrument information, the start time of the data record, and the number of samples. The data record contains the record number, 8 sample values and a checksum. This format uses a separate file for each component of each station.

ISAM-PITSA

Class: 2,4 Platform: Unix

Indexed Sequential Access Method (ISAM) is a commercial database file system designed for easy access. PITSA bases its internal file structure for digital waveform data on ISAM. This structure is often referred to as the ISAM format, but it should not be confused with the underlying database engine. An ISAM-PITSA file system consists of two database files containing the headers and the indexing information for all traces, and at least one trace file per channel. The trace file is a binary image of the floating-point data that can in principle be accessed independently. All files in an ISAM-PITSA file system have the same file name base. The extensions are ".nx0" and ".dt" for the database files, and ".001", ".002", etc. for the trace files.

Ismes

Class:1 Platform: PC

Format used by Italian Ismes recorders.

Kinometrics formats

Class:1 Platform: PC

Kinometrics have several binary formats although the two main formats are for the DataSeis recorders and the K2 class recorders.

Lennartz

Class: 1 Platform: PC

10. Seismic Data Formats, Archival and Exchange

Format for Lennartz recorders. The most common is the Mars88 format although there is also a format used with the older tape recorders.

Nanometrics

Class: 1 Platform: PC

Format used by Nanometrics recorders. The most common format is the Y-format.

NEIC ORFEUS

Class: 2 Platform: PC

The NEIC ORFEUS program SONIC1 can be used to search, display, and write data from the NEIC Earthquake Digital Data CD-ROMs (NEIC Waveform Catalog, 1991). Digital waveform data in ASCII contain two types of records: header records and data records. A header record contains station information, the start time of the data, sample rate, and the parameters of the transfer function. Data records contain the actual data retrieved from CD-ROM. Each data record is preceded by the number of data points contained in the data record. For more information, see the documentation on the NEIC ORFEUS SONIC Program Disk .

PDAS

Class: 1 Platform: PC

The format used by the Geotech PDAS recorders. This format has seen more use than just for the recorder output and there are examples of whole data sets converted to PDAS format.

PITSA BINARY

Class 2,3 Platform: PC and UNIX

In order to facilitate portability and to permit every user to write their own conversion routines without having to purchase commercial 3rd party software, a new format called BINARY has been added to PITSA's I/O. It is simply a binary image of the internal representation of data in PITSA, without the database overhead of the ISAM format. Another advantage to BINARY format is that it makes exchange of data files across platforms fairly easy. It is only necessary for the user to provide a code to do any required byte swapping. For a transitional period, fully equivalent I/O for both ISAM and BINARY routines have been implemented in both the PC and the Sun versions of PITSA, but the ISAM format will disappear eventually.

Each file consists of a short file header followed by as many data blocks as there are traces. Everything is binary. The file header consists of:

1. NCHANNELS: a long integer containing the number of channels in the file.
2. SIZE[]: An array of long integers of dimension NCHANNELS. Each element SIZE[i] contains the block size for block i, in bytes. In this context, block size of the i-th block means the size of the i-th trace header plus the size of the i-th trace.
3. BLOCK[i], for i = 1 to NCHANNELS: One block per trace. Each block consists of a binary image of the data header (as described in file data.h) followed by the binary image of the trace data.

Public Seismic Networks

Class: 1,2 Platform: PC

This format is used both as a recording and analysis format by Public Seismic Networks

SAC

Class 2 Platform: Unix

Seismic Analysis Code (SAC) is a general-purpose interactive program designed for the study of time sequential signals [SAC]. Emphasis has been placed on analysis tools used by research seismologists. A SAC data file contains a single data component recorded at a single seismic station. Each data file also contains a header record that describes the contents of that file. Certain header entries must be present (e.g., the number of data points, the file type, etc.). Others are always present for certain file types (e.g., sampling interval, start time, etc. for evenly spaced time series). Other header variables are simply informational and are not used directly by the program. Although the SAC analysis software only runs on Unix platforms and the general format is binary, there is also an ASCII version that can be used on any platform.

SEED

Class 3 Platform: All

The Standard for the Exchange of Earthquake Data (SEED) format was developed within the FDSN. The first set-up was designed at the U.S. Geological Survey's National Earthquake Information Center (NEIC) and Albuquerque Seismic Laboratory (ASL), primarily for the exchange of unprocessed waveform data. SEED was adopted by the Federation of Digital Seismographic Networks (FDSN) in 1987 as its standard. IRIS has also adopted SEED, and uses it as the principal format for its datasets. SEED uses four types of control headers:

- volume identifier headers;
- abbreviation dictionary headers;
- station headers;
- time-span headers.

Each header can use several blockettes - individual portions of information that are header specific - that conform to the organization rules of their volume type. Some blockettes vary in length and can be longer than the logical record length. Data fields in control headers are formatted in ASCII, but data fields (in data records) are primarily formatted in binary. The full description can be found in the SEED reference Manual [SEED].

It is worth pointing out that formats (such as SEED) designed to handle the requirements of international data exchange are seldom suited to the needs of individual researchers. Thus, the wide availability of software tools to convert between SEED and a full suite of Class 2 formats is crucial for its success.

A number of the present generation data acquisition systems (e.g., Quanterra, Nanometrics) produce data in SEED volumes only (miniSEED), without any of the associated control header information. Software packages have been developed to produce full SEED volumes from miniSEED volumes (e.g., SeedStuff). At the ODC, a package has recently been developed and will be distributed as a general tool.

10. Seismic Data Formats, Archival and Exchange

SEISAN

Class 2 Platform: All

The SEISAN binary format is used in the seismic analysis program SEISAN (<http://www.ifj.uib.no/seismo/software/seisan.html>). This program was developed at the Institute of Solid Earth Physics at the University of Bergen, Norway. The format consists of a main header describing all channels. Each channel then follows with a channel header with basic information including response. SEISAN can read the binary SEISAN files written on any platform. The SEISAN analysis system can also use GSE as a processing format.

SeisGram ASCII and binary

Class: 2 Platform: PC

Time series are contained in sequential, formatted ASCII files or sequential binary files. The SeisGram software (Lee, 1995) also reads fixed-record-length files using the BDSN Direct Access format. The following header information is included in both the ASCII and the binary data files:

File type, Data format, Network, Station and instrument identifier, Type of recording, Date, Event number, Orientation of the Y component, Time unit per sample, Sample rate, Amplitude units, Amplitude units per digital count, Start time, Number of samples, Comment on event and data, Time series processing history.

The ASCII files should be opened with "sequential access, formatted" format options. All header entries except start time are written with a single value on each line. The binary files are designed for compactness and fast access. Binary files should be opened with "sequential access, binary" format options. SeisGram's Direct Access data files are designed to store large sets of binary, direct access data from the BDSN (the network, not the format). The data in the file is identical to the data stream from the telemetry system, except for the addition of an eight-record header to identify uniquely the recording source, start time, and format. The Direct Access files should be opened with the "direct access, binary" format options.

Sismalp

Class: 1 Platform: PC

Sismalp is a widespread French data seismic recording system.

Sprengnether

Class: 1 Platform: PC

Format used by Sprengnether recorders.

SUDS

Class: 1,2,4 Platform: PC

SUDS stands for "The Seismic Unified Data System". The SUDS format was launched to be a more well thought out format useful for both recording and analysis and independent of any particular equipment manufacturer. The format has seen widespread use, but has lost some momentum, partly because it is not made platform-independent.

10.5 Format conversions

10.5.1 Why convert?

Ideally, we should all use the same format. Unfortunately, as the previous descriptions have shown, there are a large number of formats in use. With respect to parameter formats, one can get a long way with HYPO71, Nordic and GSE/ISF formats for which converters are available, such as in the SEISAN system. For waveform formats, the situation is much more difficult. First of all, there are many different formats and, since most are binary, there is the added complication that some will work on some computer platforms and not on others. This is a particular problem with binary files containing real numbers as for example, the SeisGram format. Additional problems are that: some formats have seen slight changes and exist in different versions; different formats have different contents so not all parameters can be transferred from one format to another; and conversion programs might not be fully tested for different combinations of data.

Many processing systems require a higher level format than the often primitive recording formats which is probably the most common reason for conversion; a similar reason is to move from one processing system to another. The SEED format has become a success for archival and data exchange, but it is not very useful for processing purposes, and almost unreadable on PCs. So it is also important to be able to move down in the hierarchy. Therefore, the main reasons for format conversion are to move:

- upwards in the hierarchy of formats for the purpose of data archiving and exchange;
- downward from the archive and exchange formats for analysis purposes;
- across the hierarchy for analysis purposes;
- from one computer platform to another.

10.5.2 Ways to convert

There are essentially two ways of converting. The first is to request data from a data center in a particular format or to log into a data center and use one of their conversion programs. The other more common way is to use a conversion program on the local computer. Such conversion programs are available both as free standing software and as part of processing systems. Equipment manufactures will often supply at least a program to convert recorder data to some ASCII format and often also to some more standard format as SUDS.

10.5.3 Conversion programs

Since conversion programs are often related to analysis programs, Tab. 10.5 lists some of the better-known analysis systems and the format they use directly.

Tab. 10.5 Examples of popular analysis programs.

Program	Author(s)	Input format(s)	Output format(s)
CDLOOK	R.Sleeman	SEED	SAC, GSE
Geotool	J.Coyne	CSS, SAC, GSE	CSS, SAC, GSE
PITSA	F.Scherbaum, J.Johnson	ISAM, SEED, Pitsa binary,	ISAM, ASCII

10. Seismic Data Formats, Archival and Exchange

		GSE, SUDS	
SAC	LLNL	SAC	SAC
SEISAN	J.Havskov, L. Ottemöller	SEISAN, GSE	SEISAN, GSE, SAC
SeismicHandler	K.Stammler	q, miniSEED, GSE, AH, ESSTF	q, GSE, miniSEED
SNAP	M.Baer	SED, GSE	SED, GSE
SUDS	P.Ward	SUDS	SUDS
Event	M.Musil	ESSTF, ASCII	ESSTF, ASCII
SeisBase	T.Fischer	ESSTF, Mars88, GSE	GSE

An overview of available format conversion programs can be found on the ORFEUS Web pages under ORFEUS Seismological Software Library (<http://orfeus.knmi.nl/wirjung.groups/wg4/index.html>). Here we present just a few packages in alphabetical order. Only those programs are mentioned which are able to read at least one of the formats mentioned in sub-Chapter 10.4.

Codeco

Program **codeco** was written by U. Kradolfer and modified by K. Stammler and K. Koch. Input files can be in SAC binary or ASCII, or GSE formats. Output formats are: integer or compressed GSE1.0 or GSE2.0, SAC binary or ASCII, and miniSEED. **Codeco** is available through the SZGRF software library (<ftp://ftp.szgrf.bgr.de/pub/software>).

Convseis

Converts 14 data formats on PCs like GSE1.0 and GSE2.0 INT, PCEQ, SEG Y and SUDS. **Convseis** has been written by L. Oncescu and M. Rizescu.

isam2gse

Data in ISAM format can be converted to GSE format by using the program **isam2gse**. The code is available through the SZGRF software library (<ftp://ftp.szgrf.bgr.de/pub/software>).

ESSTF to GSE

Program **len2gse2**, written by B. Ruzek (Geophysical Institute, Prague) converts multiplexed ESSTF binary format, Mars88 binary format or ASL ASCII format in *data_file* to the GSE2.0 CM6 compression format. The user can select the time window and mask channels and streams. The code is written in C++.

GSE to SEED

Program **gse2seed**, developed by R. Sleeman (Orfeus Data Centre, de Bilt), converts a GSE2.X file to the SEED2.3 format. Multiple traces are handled. For each WID2 section, the GSE file must contain corresponding data types STATION, CHANNEL and RESPONSE.

PASSCAL package

The PASSCAL package was written by P. Friberg, S. Hellman, and J. Webber, developed on SUN under SunOs4.1.4, compiled under Solaris 2.4 and higher and also under LINUX. It converts RefTek to SEGY and miniSEED. Program **pql** provides a quick and easy way to view SEGY, SAC, miniSEED or AH seismic data. **pql** operates in the X11 window environment. The package is available from the PASSCAL instrument center (<http://www.passcal.nmt.edu>) at New Mexico Tech., Socorro.

Preproc

Preproc has been designed to assist the seismologist who wishes to analyze large sets of raw digital data that need to be preprocessed in some standard way prior to the analysis. Preproc was written by Miroslav Zmeskal for the ISOP project in the period 1991-1993. It was rewritten recently in the object-oriented form. As a by-product, **preproc** can perform data conversion from GSE / PITSA ISAM to GSE / PITSA ISAM. In the near future new input/output formats will be implemented (ESSTF, miniSEED). **preproc** was successfully compiled on HP, SUN, Linux and DOS. Program package **preproc** and a detailed Manual are available through the ORFEUS Seismological Software Library

Rdseed

Rdseed reads from the input tape or file in the SEED format. According to the command line function option specified by the user, **rdseed** will read the volume and recover the volume table of contents (-c), the set of abbreviation dictionaries (-a), or station and channel information and instrument response table (-s). In order to extract data from the SEED volume for analysis by other packages, the user must run **rdseed** in user prompt mode (without any command line options). As data is extracted from the SEED volume, **rdseed** looks at the orientation and sensitivity of each channel and corrects the header information on request. Implemented output formats are (option d): SAC, AH, CSS 3.0, miniSEED and SEED. A Java version of rdseed is to be released in 2001. **Rdseed** was developed by Dennis O'Neill and Allen Nance, IRIS DMC.

SeedStuff

SeedStuff is a set of basic programs provided by the GEOFON DMS software library in Potsdam (<ftp://ftp.gfz-potsdam.de/pub/home/st/GEOFON/software>) to process and compile raw data from Quanterra, Comserv and RefTek data loggers. The goal is to check and extract data from station files/tapes to miniSEED files and to assemble miniSEED files to full SEED volumes. The SeedStuff package was written by Winfried Hanka and compiled on the SUN, HP and Linux. The following tools are available:

- extr_qic:** extracts multiplexed raw Quanterra station tapes to demultiplexed miniSEED files containing only one station / stream / component;
- extr_file:** like extr_qic for multiplexed miniSEED, RefTec files;
- extr_fseed:** disassemble full SEED tapes. SEED headers are skipped, data are stored into station / stream / component files;
- check_seed:** checks the contents of miniSEED data files or tapes ;

10. Seismic Data Formats, Archival and Exchange

check_qic: analysis the contents of a Quanterra data tape;

copy_seed: assembles a full SEED volumes from miniSEED files for a given set of station / stream / component defined in the copy_seed.cfg configuration file

make_dlsv: generates a dataless (header only) SEED volume for a set of station/stream/component defined in copy_seed.cfg.

SEED to GSE

There is no special program developed for converting either full SEED volumes or miniSEED files to the GSE format. Such a package would be strongly needed for providing data in the GSE format by the AutoDRM services.

On the SUN platform, program CDLOOK (see 11.5.2.2) can read full SEED volumes and write traces in the GSE format. This program can be downloaded from <ftp://orfeus.knmi.nl/pub/software>.

SEISAN

The SEISAN analysis system has about 40 conversion programs, mostly from some binary format to SEISAN. The SEISAN format can then be converted to any standard format like SEED, SAC or GSE. SEISAN has format converters for most recorders on the market including Kinometrics, Nanometrics, Teledyne, GeoSig, Reftek, Lennartz, Güralp and Sprengnether.

Acknowledgments

The authors acknowledge with thanks the careful review by Bruce Presgrave of the US Geological Survey. It has improved both the language of the original draft and provided useful references to the Earthworm system. Thanks go also to Xiaoping Yang who kindly provided the links to the data bases of the Center for Monitoring Research and the CTBTO.

Special references

- [CSS] Anderson, J., W. Farrell, K. Garcia, J. Given, and H. Swanger, Center for Seismic Studies Version 3 Database: Schema Reference Manual, SAIC Technical Report C90-01, 1990.
- [IDC3.4.1] Formats and Protocols for Messages, Rev. 3, 2001.
- [GSE] Provisional GSE2.1 Message Formats & Protocols, 1997. Operations Annex 3, GSETT-3.
- [LEN] SAS-58000 User's Guide and Reference Manual, 1986. Lennartz electronic GmbH
- [SAC] W.C. Tapley & J.E. Tull, 1992. SAC - Seismic Analysis Code. LLNL, Regents of the University of California
- [SEED] Standard for the Exchange of Earthquake Data, 1992. Reference Manual, SEEDFormat v2.3, FDSN, IRIS, USGS