

1 Draft advice to the European Medicines Agency from the 2 clinical trial advisory group on Clinical trial data formats 3

4 The clinical trial advisory group at their TC on 7th March 2013 provided
5 advice as follows:

6 **1. The following definitions were agreed**

7 1.1 This advice refers to all data recorded in a clinical trial (as part of documents and data
8 aggregated or at patient level) that can be stored electronically and associated documentation
9 (additional information that identifies and characterises the data properties such as dataset keys,
10 variable definition, terminology, code lists) that is submitted by Applicants to the Agency throughout
11 the life-cycle of medicinal products. The policy will be applied prospectively for future submissions to
12 the Agency; it may include old clinical trial data.

13 1.2 Data formats, in this advice, refer to the organisation of information according to pre-set
14 specifications that facilitate the storage, exchange, access, comprehension, analysis and archive of
15 clinical data. It includes both the type of electronic files and the structure of the files, as well as
16 associated documentation.

17 The data and associated documentation concerned by this policy may or may not be sourced via
18 electronic tools (e.g., paper or electronic case report forms), and are subsequently submitted and
19 stored electronically.

20 All statements in this advice are made in the consideration that CTAG1 rules for patient data
21 confidentiality and anonymisation are applied and effective, and that CTAG5 legal rules are strictly
22 followed. As a consequence, there will be no more reference to the CTAG1 and CTAG5 rules in the rest
23 of this advice.

24 **2. There is a need to define data formats**

25 The choice of formats should neither imply delays in the information to be made available nor impose
26 unnecessary burden to the stakeholders.

27 Formats may be different depending on the type of information to be made publicly available and the
28 intended use of it, and data should be made available irrespectively.

29 As there are not universally agreed formats, a minimum set of rules should be defined, including:

- 30 • An indexed list of all clinical trials present in the submissions shall be provided (if not already
31 available in the table of content of the submission dossier) so the data of the overall clinical
32 program is tracked. In this list, clinical trials should ideally be identified by a unique trial
33 identifier. This identifier could be either: the EudraCT number (but it would only cover trials
34 conducted in Europe and it is not commonly referenced in journals and published articles), the
35 NIH clinicaltrials.gov website registry number (commonly referred to in the literature), the
36 ISRCTN registry number or a number provided by the applicant at the time of submission. It is
37 thought to be useful to be able to link back clinical trials to journal article information.

- 38 • Data shall be published in the format they have been submitted and evaluated and no
39 conversion of formats will be done by either the marketing authorisation holder or the
40 European Medicines Agency (EMA).
- 41 • Consistency of formats throughout the life cycle of the medicinal products is not mandatory but
42 should be sought when achievable, e.g. for contemporaneous studies.
- 43 • Documents containing data should be human readable and searchable by anyone requesting
44 the data from the EMA.
- 45 • There was a request that analysis of patient-level data could be done in Excel.
- 46 • Patient-level data should be accompanied by associated documentation that allows quickly
47 grasping the data and processing it. This documentation, which includes metadata (=
48 'structured data about data'), should ideally be machine readable. For example, the
49 documentation explains the structure of the data (e.g., what information is contained in each
50 dataset), gives the definition of data elements (e.g., '1' corresponds to 'male' and '2' to
51 'female'), and provides the context to interpret correctly the data, to allow further analyses,
52 without needing for additional information from neither the marketing authorisation holder nor
53 the EMA.
- 54 • Formats should be chosen so that data is readable with open source, non-proprietary software
55 (but not necessarily free): that includes, but is not limited to, portable document format (PDF)
56 for text documents such as clinical study reports, SAS transport file format (XPT) for datasets
57 and programs (as opposed to SAS format which is proprietary), and extensible markup
58 language (XML) format for associated documentation on data. It would be easier for Industry
59 in general if these requirements are the same as FDA's, although it is not favoured by small-
60 and medium-sized enterprise if at a non-negligible cost.

61 **3. What is to be included in data formats**

62 Assuming that data privacy protection has been ensured for all data made available publicly, cCertain
63 information such as CT scans, MRI and other imaging, interviews, genetic/genomic data can bring
64 useful information and should be in the scope of discussion for data formats. However, that particular
65 type of data is contained in large files; thus its transport, storage and access might cause serious
66 informatics problems.

67 Three levels of clinical trial information, data and associated documentation shall be included.

- 68 • Level 1: for each product, a full list of clinical trials, including a unique study identifier; these
69 lists should be fully searchable and could be connected to the European Public Assessment
70 Reports. This is separate to information stored in the EUdraCT database.
- 71 • Level 2: for each study, full clinical study report (CSR) according to ICH E3, including all
72 appendices, as detailed in ICH E3 (study information, patient data listings and case report
73 forms [CRF]).
- 74 • Level 3: for each study, individual patient data sets (including individual patient data) and
75 additional results used for the evaluation of the drug (if not covered by Level 2),
76 documentation explaining the structure and content of datasets (e.g., annotated CRF, variable
77 definitions, data derivation specifications, dataset define file), test outputs, SAS logs and SAS
78 programs

79 Elements included in the three levels of data listed above may need to be modified in special
80 circumstances driven by confidentiality or legal aspects.

81 **4. Formats recommended**

82 In general, to avoid delays any format shall be acceptable for all data until the policy is applied by
83 stakeholders. The data shall be published in the format they are available at present.

84 In terms of the different types of data described in the previous section, Level 1 data should be
85 searchable. PDF is recommended.

86 For Level 2 data (CSR and appendices, according to ICH E3), it should also be searchable. PDF is
87 recommended. Of note, old CSRs may not fully comply with the current ICH E3 format. In this case, it
88 will be acceptable to provide the CSR in the original format in which it was written.

89 Individual patient data and associated documentation (Level 3) shall be published in the format they
90 are available at time of submission. That can be according to CDISC standards, and there was general
91 agreement that Applicants will move progressively to an increase use of CDISC standards.

92 It was recognised that CDISC have defined useful formats: SDTM for data tabulations, ADaM for
93 analysis datasets, and define xml for metadata. The recommendation is for all these to be submitted to
94 the Agency, but not ODM, which is a transport format for data management. SDTM-annotated CRF
95 would also be very useful for data re-analysis. It was acknowledged that CDISC implementation guides
96 can be interpreted in different ways by Applicants, therefore EMA should define clear requirements in
97 relation to these guides.

98 If other formats can be used, EMA should define minimal requirements of more basic formats, such as
99 the following: clinical data should be submitted in rectangular tables, in a comma-separated values
100 (CSV) format; associated metadata should contain at least one table with all datasets, all variables and
101 their meanings, also possibly associated code lists, and another table with all codes and decodes, and
102 the variables they relate to.

103 Individual data such as CRF data in PDF format are not useful as they will require substantial
104 manpower for reloading in another usable format). However, PDF scans of printed out CRFs might be
105 the minimal standard which is realisable even in a small academic institution or a small- and medium-
106 sized enterprise, in order not to add unnecessary financial and resource burden to the marketing
107 authorisation holder. The general view is that re-formatting of old data should be not requested by
108 EMA; however, some are of the opinion that EMA should ask the marketing authorisation holder to
109 provide the data in a format which is machine-readable and can be done with a non-proprietary
110 software.

111 Harmonisation of formats such as CDISC SDTM and ADAM is of course desirable as this expands the
112 usefulness of the data made available. This exercise shall be progressively implemented in a
113 collaborative way between CDISC and EMA to ensure consistency and versioning control.

114 Sustainability of a chosen standard might also require reducing the speed of versioning and ensuring
115 availability of software adapted to the subsequent changes of the formats. EMA guidance on formats
116 may not follow the evolution of CDISC modifications at the same rhythm if it imposes too much burden
117 on applicants. This will reduce the potential for re-formatting should a newer version be required.

118 Formats used across a number of studies for the same product do not need be compatible, although it
119 will be a bonus when it can be achieved. For the datasets there is a need to:

120 • Harmonise a reference format worldwide

121 • Maintain versioning over time

122 A point to discuss further concerns mixed formats acceptability, e.g. for fixed combination of old and
123 new active substances or hybrid mixed submission, when both clinical data from old studies and from
124 new clinical trials are included.

125 **5. Who should adhere to the agreed formats**

126 The formats agreed are to be adhered to by all stakeholders and also for locally run clinical trials
127 outside Europe if they become part of a submission to EMA. The Applicants should ensure correct
128 implementation of the formats and should also consider implication of terms translations from different
129 languages.

130 For clinical trials owned in different measure by multiple partners (e.g. public-private partnerships),
131 the above points should be taken into account from the beginning of the clinical studies. This concerns
132 data that are part of studies that are submitted to the Agency and where the marketing authorisation
133 holder is legally permitted to share the data.

134 **6. Timelines for format implementation**

135 • While it seems reasonable to gain experience with formats of individual patient data (Level 3),
136 it is not recommended to have a test period for clinical study reports, because the format of
137 the CSRs, i.e. ICH E3, is in effect since 1996. Therefore the format for CSRs (Level 2) - and for
138 Level 1 - can be mandatory from the implementation of the policy.

139 • Pro-active adoption of standard formats for Level 3 data: as this has to be mandatory for the
140 sake of fairness and clarity for all stakeholders, it is advised to start gradually to acquire
141 experience and then mandate formats after a trial period for all new studies submitted.

142 • At the end of this trial period, all levels of data can be released at the same time.

143 **7. International harmonisation across regulatory agencies**

144 The EMA is leading in terms of policy but global alignment and harmonisation are critical steps in the
145 future process. A global consultation of formats is recommended at the ICH level (for human products
146 and at the VICH level for veterinary products). The list of elements discussed in Section 3 and the
147 corresponding formats discussed in Section 4 need to be included in that consultation. Communication
148 with other national medicines agencies would also be beneficial. The policy should also aim at
149 implementing what will be widely used in future to further standardise the process and prevent any re-
150 formatting.

151 Under e-CTD, PDF, XML and other standards are allowed in MAA. ISO, CEN and CDISC to define CSRs
152 harmonised standards.

153 **8. References**

154 ICH E3 Structure and Content of Clinical Study Reports

155 http://www.ich.org/fileadmin/Public_Web_Site/ICH_Products/Guidelines/Efficacy/E3/E3_Guideline.pdf

156 Clinical Data Interchange Standards Committee

157 <http://www.cdisc.org>

30 April 2013

Advice to the European Medicines Agency from the clinical trial advisory group on Clinical trial data formats (CTAG2) – meeting 3 outcome with comments and amendments

Annex I - Comments from participants below may or may not have been made on behalf of the organisation they are affiliated with.

Line number	Comment and Changes proposed	Name	Affiliation
61	<p>Comment: The title could be made clearer by making reference to "types" or "elements" of data, as this is the key section where the specific types of data are specified and divided into three levels (Level 1, 2, 3).</p> <p>Proposed change (if any): Rename section as "What types of data (and in what format) are to be included".</p>	Peter Doshi	Johns Hopkins University School of Medicine

Line number	Comment and Changes proposed	Name	Affiliation
61	<p>Comment: The clinical study report (CSR, defined by ICH E3) and SAS datasets are two important types of trial data, but they are not the only types of trial data. As a major function of CTAG2 is to define the types of clinical trial data that--if suitably de-identified--can be prospectively released, it is important to include within the CTAG2 advice document a clear table that outlines the many types of data relevant to understanding and interpreting clinical trials.</p> <p>For each type of data that exists, EMA may already routinely be requesting these data from marketing authorization applicants, or they may consider requesting these data in the future. Whether EMA requests or intends to routinely request these data will clearly impact the ability of EMA to prospectively release these data. Those writing the EMA's draft policy should know whether EMA plans to hold such information itself.</p> <p>Part of the feasibility of prospective data release also depends on the format of the data (e.g. paper, PDF, or proprietary computer file format). Those writing the EMA's draft policy should have this information in order to better decide the feasibility and mechanism by which the data may be de-identified and released.</p> <p>Different de-identification issues may arise with different types of data. For instance, a trial's statistical analysis plan (ICH E3 section 16.1.9) has no information about patients, and presumably does not require de-identification before prospective release, whereas the list of study investigators (ICH E3 section 16.1.4) and a patient-level SAS dataset may require de-identification prior to release. There is a need to record the need for CTAG1 advice on de-identification for the types of data that CTAG2 identifies.</p>	Peter Doshi	Johns Hopkins University School of Medicine

Line number	Comment and Changes proposed	Name	Affiliation
	<p>Proposed change (if any): I propose adding a table to section 3 titled "Table of types of data" to this section that consists of 6 columns and accomplishes the above described needs.</p> <p>The columns would be as follows: Column 1: Type of data Column 2: Routinely requested by EMA (yes/no)? Column 3: Format (e.g. paper, PDF, electronic dataset such as SAS XPORT, etc.) Column 4: Level of data (i.e. Level 1, 2, or 3) Column 5: Need for CTAG1 advice on risk of participant re-identification (yes/no)? Column 6: CTAG1's advice on who, when (pre-submission or post-submission to EMA), and how to ensure acceptably low risk of re-identification (leave blank until CTAG1's advice is received)</p> <p>===== For Column 1 (type of data), I suggest it be populated with the following:</p> <ul style="list-style-type: none"> - For each product, a full list of clinical trials, including a unique study identifier, the study title, the interventions and the indication studies; these lists should be fully searchable. The studies should be connected to the European Public Assessment Reports, to EUDRACT and to clinicaltrials.gov. [LEVEL 1] - Full Clinical Study Report (CSR) including all appendices defined by ICH E3 (http://www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2009/09/WC500002832.pdf) [LEVEL 2] <ul style="list-style-type: none"> - ICH E3 1-15 Core report - ICH E3 16.1 Study information - ICH E3 16.1.1 Protocol and protocol amendments 	Peter Doshi	Johns Hopkins University School of Medicine

Line number	Comment and Changes proposed	Name	Affiliation
	<ul style="list-style-type: none"> - ICH E3 16.1.2 Sample case report form (unique pages only) - ICH E3 16.1.3 List of IECs or IRBs (plus the name of the committee Chair if required by the regulatory authority) - Representative written information for patient and sample consent forms - ICH E3 16.1.4 List and description of investigators and other important participants in the study, including brief (1 page) CVs or equivalent summaries of training and experience relevant to the performance of the clinical study - ICH E3 16.1.5 Signatures of principal or coordinating investigator(s) or sponsor's responsible medical officer, depending on the regulatory authority's requirement - ICH E3 16.1.6 Listing of patients receiving test drug(s)/investigational product(s) from specific batches, where more than one batch was used - ICH E3 16.1.7 Randomisation scheme and codes (patient identification and treatment assigned) - ICH E3 16.1.8 Audit certificates (if available) - ICH E3 16.1.9 Documentation of statistical methods - ICH E3 16.1.10 Documentation of inter-laboratory standardisation methods and quality assurance procedures if used - ICH E3 16.1.11 Publications based on the study - ICH E3 16.1.12 Important publications referenced in the report - ICH E3 16.2 Patient data listings - ICH E3 16.3 Case Report Forms - Other documents containing contextual information not identified in ICH E3: <ul style="list-style-type: none"> - certificate(s) of analysis [LEVEL 2] - investigator's brochure [LEVEL 2] - manual of operations and procedures [LEVEL 2] - annotated CRFs [LEVEL 3] 		

Line number	Comment and Changes proposed	Name	Affiliation
	<ul style="list-style-type: none"> - Electronic database of Individual Participant Data (IPD) <ul style="list-style-type: none"> - patient-level dataset (raw and derived) [LEVEL 3] - analysis datasets [LEVEL 3] - Contextual information to understand and work with electronic database of IPD <ul style="list-style-type: none"> - dataset specifications (metadata which describes the variable labels, variable descriptions, code lists and formats) - SAS programs [LEVEL 3] - SAS logs [LEVEL 3] - test outputs [LEVEL 2] - Original participant level records <ul style="list-style-type: none"> - filled out (completed) CRFs for all trial participants [LEVEL 3] - laboratory reports for all trial participants [LEVEL 3] - medical records and diagnostic reports for all trial participants obtained as part of trial procedures [LEVEL 3] - Documents related to clinical trials often created by trial sponsors <ul style="list-style-type: none"> - Marketing Assessments - Email correspondence - Meeting minutes - Records of the Data Monitoring Committee (also known as DSMB) e.g. adjudication committee 		

Line number	Comment and Changes proposed	Name	Affiliation
	<ul style="list-style-type: none"> - ICH E3 16.1.1 Protocol and protocol amendments - ICH E3 16.1.2 Sample case report form (unique pages only) - ICH E3 16.1.3 List of IECs or IRBs (plus the name of the committee Chair if required by the regulatory authority) - Representative written information for patient and sample consent forms - ICH E3 16.1.4 List and description of investigators and other important participants in the study, including brief (1 page) CVs or equivalent summaries of training and experience relevant to the performance of the clinical study - ICH E3 16.1.5 Signatures of principal or coordinating investigator(s) or sponsor's responsible medical officer, depending on the regulatory authority's requirement - ICH E3 16.1.6 Listing of patients receiving test drug(s)/investigational product(s) from specific batches, where more than one batch was used - ICH E3 16.1.7 Randomisation scheme and codes (patient identification and treatment assigned) - ICH E3 16.1.8 Audit certificates (if available) - ICH E3 16.1.9 Documentation of statistical methods - ICH E3 16.1.10 Documentation of inter-laboratory standardisation methods and quality assurance procedures if used - ICH E3 16.1.11 Publications based on the study - ICH E3 16.1.12 Important publications referenced in the report - ICH E3 16.2 Patient data listings - ICH E3 16.3 Case Report Forms 	Peter Doshi	Johns Hopkins University School of Medicine

Line number	Comment and Changes proposed	Name	Affiliation
	<ul style="list-style-type: none"> - Other documents containing contextual information not identified in ICH E3: <ul style="list-style-type: none"> - certificate(s) of analysis [LEVEL 2] - investigator's brochure [LEVEL 2] - manual of operations and procedures [LEVEL 2] - annotated CRFs [LEVEL 3] - Electronic database of Individual Participant Data (IPD) <ul style="list-style-type: none"> - patient-level dataset (raw and derived) [LEVEL 3] - analysis datasets [LEVEL 3] - Contextual information to understand and work with electronic database of IPD <ul style="list-style-type: none"> - dataset specifications (metadata which describes the variable labels, variable descriptions, code lists and formats) - SAS programs [LEVEL 3] - SAS logs [LEVEL 3] - test outputs [LEVEL 2] - Original participant level records <ul style="list-style-type: none"> - filled out (completed) CRFs for all trial participants [LEVEL 3] - laboratory reports for all trial participants [LEVEL 3] - medical records and diagnostic reports for all trial participants obtained as part of trial procedures [LEVEL 3] - Documents related to clinical trials often created by trial sponsors <ul style="list-style-type: none"> - Marketing Assessments - Email correspondence - Meeting minutes - Records of the Data Monitoring Committee (also known as DSMB) <p>e.g. adjudication committee</p>	Peter Doshi	Johns Hopkins University School of Medicine

30 April 2013

Advice to the European Medicines Agency from the clinical trial advisory group on Clinical trial data formats (CTAG2) – meeting 3 outcome with comments and amendments

Line number	Comment and Changes proposed	Name	Affiliation
	- Records of the Data Monitoring Committee (also known as DSMB) e.g. adjudication committee	Peter Doshi	Johns Hopkins University School of Medicine