



# 1 Best Practice

---

## 2 **Subject**

3 DDI 3.0 URNs and Entity Resolution (2009-03-21)

## 4 **Document identifier:**

5 DDIBestPractices\_URNsAndEntityResolution.doc.PDF

## 6 **Location:**

7 [http://www.ddialliance.org/bp/DDIBestPractices\\_URNsAndEntityResolution.doc.PDF](http://www.ddialliance.org/bp/DDIBestPractices_URNsAndEntityResolution.doc.PDF)

## 8 **Authors:**

9 Nikos Askitas, Janet Eisenhauer, Arofan Gregory, Rob Grim, Pascal Heus, Maarten  
10 Hoogerwerf, Wolfgang Zenk-Möltgen

## 11 **Editors:**

12 Nikos Askitas, Maarten Hoogerwerf

## 13 **Abstract:**

14 This document is not quite a best practice but rather a recommendation about  
15 appropriate architecture for the effective resolution of DDI URNs<sup>1</sup>. The  
16 recommended architecture is based on standard and tested technologies put  
17 together in order to facilitate URN resolution needs. Along the way we describe the  
18 consequences for the various parties involved and the relationship of DDI URNs to  
19 other resolution mechanisms. This document does not deal with the latter in depth,  
20 as that would be the subject of a white paper in its own right.

## 21 **Status:**

22 This document is updated periodically on no particular schedule. Send comments to  
23 editor: [ddi-bp-editors@icpsr.umich.edu](mailto:ddi-bp-editors@icpsr.umich.edu)

---

<sup>1</sup> All URNs in this document are assumed to be DDI 3.0 URNs. Unless otherwise stated, we simply say “URN” and mean “DDI 3.0 URN”.



24

25 **INTRODUCTION..... 3**

26 **1.1 Problem statement ..... 3**

27 **1.2 Terminology ..... 3**

28 **2 BEST PRACTICE SOLUTION ..... 4**

29 **2.1 Definitions ..... 4**

30 **2.2 Best Practice behavior ..... 4**

31 **2.3 Implications ..... 7**

32 **2.4 Validity ..... 7**

33 **2.5 Discussion ..... 7**

34 **3 REFERENCES ..... 10**

35 **3.1 Normative ..... 10**

36 **APPENDIX A. ACKNOWLEDGMENTS ..... 11**

37 **APPENDIX B. REVISION HISTORY ..... 13**

38 **APPENDIX C. LEGAL NOTICES ..... 14**



39

## 40 **Introduction**

41 This best practice lays out an architecture and solution to the issue of providing unique  
42 persistent identifiers for DDI entities down to the variable level.

### 43 **1.1 Problem statement**

44 A key feature of DDI 3.0 is the use of URN-based references for the purpose of identification,  
45 discovery, and reusability. To enable this to the right granularity, many URNs need to be  
46 assigned. This implies scale, mass, and volume issues that need to be addressed and solved.  
47 To answer these issues effectively, an architectural infrastructure must have the following core  
48 properties:

- 49 1. High availability, reliability, durability
- 50 2. Scalability, manageability on all levels (DDI Alliance, institution, etc.)
- 51 3. Non-invasive, low barrier, compatibility, and interoperability
- 52 4. Reliance on standard and tested general purpose technologies
- 53 5. Community-based (archival, library, academic, quantitative data)

### 54 **1.2 Terminology**

55 The key words *must*, *must not*, *required*, *shall*, *shall not*, *should*, *should not*, *recommended*,  
56 *may*, and *optional* in this document are to be interpreted as described in **[RFC2119]**. Additional  
57 DDI standard terminology and definitions are found in <http://www.ddialliance.org/definitions/>.



58

## 59 **2 Best Practice Solution**

### 60 **2.1 Definitions**

61 DNS: The Domain Name System (DNS) translates Internet domain and host names to IP  
62 addresses. It translates domain names meaningful to humans into the numerical (binary)  
63 identifiers associated with networking equipment for the purpose of locating and addressing  
64 these devices world-wide. An often used analogy to explain the Domain Name System is that it  
65 serves as the "phone book" for the Internet by translating human-friendly computer hostnames  
66 into IP addresses. For example, www.example.com translates to 208.77.188.166.

67 HTTP: Short for HyperText Transfer Protocol, the underlying protocol used by the World Wide  
68 Web. HTTP defines how messages are formatted and transmitted, and what actions Web  
69 servers and browsers should take in response to various commands. For example, when you  
70 enter a URL in your browser, this actually sends an HTTP command to the Web server directing  
71 it to fetch and transmit the requested Web page.

72 HTTPS: HTTPS stands for HyperText Transfer Protocol over SSL (Secure Socket Layer). It is a  
73 TCP/IP protocol used by Web servers to transfer and display Web content securely. The data  
74 transferred are encrypted so that they cannot be read by anyone except the recipient.

75 IP: An Internet Protocol (IP) address is a numerical identification (logical address) that is  
76 assigned to devices participating in a computer network utilizing the Internet Protocol for  
77 communication between its nodes. Although IP addresses are stored as binary numbers, they  
78 are usually displayed in human-readable notations, such as 208.77.188.166 (for IPv4).

79 URL: A URL (Uniform Resource Locator, previously Universal Resource Locator) is the unique  
80 address for a file that is accessible on the Internet. A common way to get to a Web site is to  
81 enter the URL of its home page file in a Web browser's address line. However, any file within  
82 that Web site can also be specified with a URL.

83 URN: A URN (Uniform Resource Name) is an Internet resource with a name that, unlike a URL,  
84 has persistent significance -- that is, the owner of the URN can expect that someone else (or a  
85 program) will always be able to find the resource.

86 TTL: Short for Time to Live, a field in the Internet Protocol (IP) that specifies how many more  
87 hops a packet can travel before being discarded or returned.

88

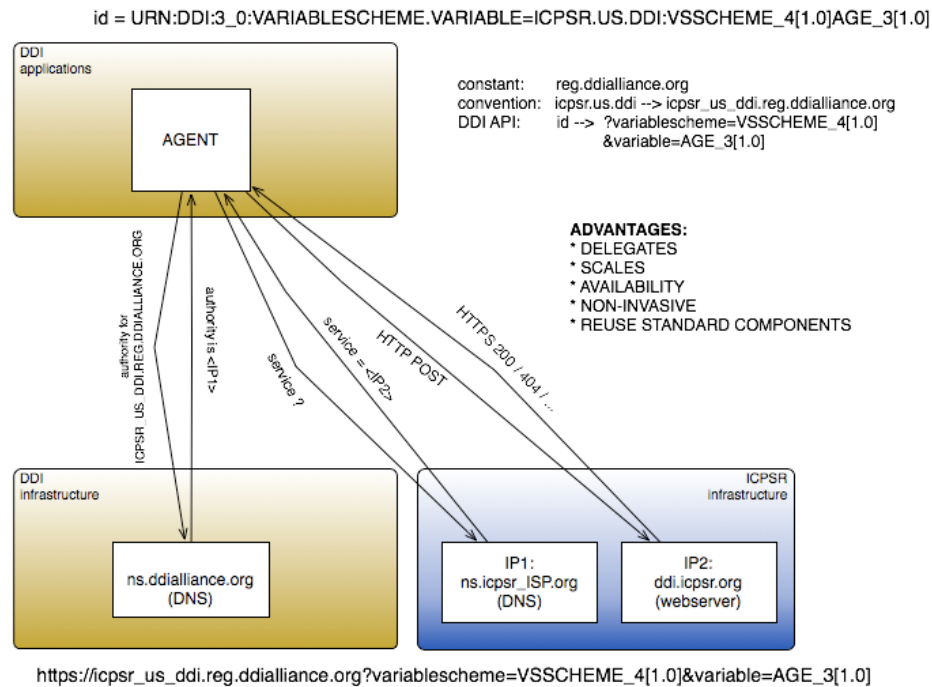
### 89 **2.2 Best Practice behavior**

90 **General description**

91 The proposed solution will be used for resolution of a great number of DDI URNs<sup>2</sup>. Combined  
 92 with the requirements mentioned in paragraph 1.1 “Problem statement,” this means delegation  
 93 of both responsibility and resolution. An existing system that implements such delegation and is  
 94 well proven is DNS. Rather than rebuilding or copying DNS, the proposed infrastructure will use  
 95 it to meet the requirements of the solution.

96 In addition to DNS, the proposed architecture will draw upon highly used techniques such as  
 97 (secured) HTTP, URLs, and query strings. Reuse of this technology acknowledges the concepts  
 98 of DDI.

99 The solution is explained here using a schematic overview. This overview demonstrates both  
 100 the components of the solution and the resolution process. The overview uses ICPSR as an  
 101 example, but in general one can substitute any DDI Agency wherever “ICPSR” appears.



102

103

**Figure 1: Overview of DNS-based resolution**

104 **Description of required infrastructure**

105 Figure 1 shows an infrastructure consisting of three components described below:

<sup>2</sup> An estimation for the amount of resolvable URNs is 1.500.000.000 (500 agencies each having 500 studies that each contain 2000 variables that each have 3 referred schemes)



## Data Documentation Initiative

- 106 • DDI applications. A typical situation is that a service or application wants to retrieve the  
107 referred material to show it to an end user. All DDI aware applications should have a  
108 central agent that takes care of resolution of the URNs.
  
- 109 • DDI Infrastructure. The DDI Alliance will need control over a name server and a domain  
110 name where it can register various sub domains for each registered DDI Agency:  
111 <agency\_id.reg>.ddialliance.org. Redundancy or load balancing of this  
112 infrastructure can require agreements with certain agencies.<sup>3</sup>
  
- 113 • Agency Infrastructure. Each agency needs to be able to register some DNS records in  
114 the name server of their ISP. In addition, they need their DDI-exposing application to  
115 implement an API that specifies the requests for certain concepts, studies, variables,  
116 questions, versions, etc. The API should also specify functionality like authentication and  
117 a TTL mechanism.

### 118 **Description of resolution process**

119 Assume that the agent from a researcher's client-application needs to resolve the following

120 URN: URN:DDI:3\_0:VARIABLESCHEME.VARIABLE=ICPSR.US.DDI<sup>4</sup>:VSSCHEME\_4[1.0]AGE\_3[1.0]

121 The following steps will result in the researcher's application's agent retrieving the metadata  
122 identified by this URN .

- 123 1. The DDI application recognizes ICPSR.US.DDI as the DDI Agency and queries DNS for  
124 ICPSR\_US\_DDI.REG.DDIALLIANCE.ORG, where it can resolve this URN.
  
- 125 2. The name server of the DDI Alliance responds with the name of the name server of  
126 ICPSR, which is NS.ICPSR.UMICH.EDU.
  
- 127 3. The name server of ICPSR maintains a list of services and the locations of these. The  
128 agent DDI application now requests the DDI resolution service and receives one or more  
129 IP address of the machines that host this service.
  
- 130 4. The agent now translates the URN into a query string that will request the identified data  
131 from the ICPSR DDI Registry. In a nutshell, this transformation will result in  
132 "?variablescheme=VSSCHEME\_4[1.0]&variable=AGE\_[1.0]". The details about this  
133 transformation can be found in section 2.5 "Discussion".

---

<sup>3</sup> See Discussion.

<sup>4</sup> The DDI Alliance is working with the community and its partners to establish a global registry of agencies that will be publicly available online -- see <http://tools.ddialliance.org/?lvl1=community&lvl2=agencyid>. It is crucial for early adopters of DDI 3.0 to immediately start using unique global identifiers. Therefore, there is now an unofficial pre-registration form that will allow organizations to request and reserve an identifier in the future registry. The conventions in the registry are of the form `icpsr.us.ddi`, `odesi.ca.ddi`, `gesis.de.ddi`, etc.



## Data Documentation Initiative

- 134 5. The agent makes a request of the DDI registry for the identified material using a HTTP  
135 or a HTTPS POST command like:  
136 [http://<ddisvc>.icpsr\\_us\\_ddi.reg.ddialliance.org/?variablescheme=  
138 VSSCHEME\\_4\[1.0\]&variable=AGE\\_\[1.0\]](http://<ddisvc>.icpsr_us_ddi.reg.ddialliance.org/?variablescheme=<br/>137 VSSCHEME_4[1.0]&variable=AGE_[1.0])

139 Additional parameters can be required: authentication to allow authorization, time-to-live  
140 to prevent circular references, etc.

- 141 6. The DDI registry uses the request parameters and determines whether the requested  
142 material is available and whether the requesting agent is authorized to view it. It will  
143 send its response via HTTP(S) and contain HTTP response codes and DDI XML  
144 (possibly contained within an XML container to facilitate additional parameters).

### 145 **2.3 Implications**

- 146 • The ddialliance.org creates a subdomain called reg.ddialliance.org, which maintains a  
147 domain name for each registered Agency (this implies modifying the current tentative  
148 registry appropriately). For example, this means that the current entry for ICPSR  
149 (icpsr.us.ddi) should result in the subdomain icpsr\_us\_ddi.reg.ddialliance.org to be  
150 registered as a record `dns*.ddialliance.org` in the name servers of the DDI  
151 Alliance.
- 152 • Each DDI provider registers their DDI service(s) in their own (or their Internet provider's)  
153 DNS servers. For example ICPSR needs to register at least one SRV record in their  
154 name server NS.ICPSR.UMICH.EDU
- 155 • The DDI Alliance should encourage the development of reusable libraries that take care  
156 of the resolution and/or implement the API for querying DDI providers. This will improve  
157 the quality of the resolution and lower the barriers for adopting the DDI 3.0 standard.

### 158 **2.4 Validity**

159 The choice for DNS allows the use of highly standardized protocols and software. Its reliability is  
160 proven by the current Internet infrastructure. It allows redundancy and scalability. It avoids the  
161 need to build software and can reuse existing hardware. There is little additional administration  
162 required (since the DDI already maintains a registry).

### 163 **2.5 Discussion**

#### 164 **Agency identifiers**

165 Agency identifiers can be reused to register subdomain records in DNS. Currently, the Agency  
166 identifiers are structured like ICPSR.US.DDI. This raises two issues:

- 167 • Do the dots imply additional hierarchy within DNS? Replacement of the dot (.) by an  
168 underscore (\_) can help to overcome this problem, either as a system to register an



## Data Documentation Initiative

169 agency identifier or only as a translation to determine the correct subdomain name:  
170 ICPSR\_US\_DDI.

- 171
- How unique are these identifiers? To be complete, these identifiers should be integrated  
172 within a standardized domain. This can be done by prefixing them with, for example,  
173 URN:DDI:AGENCY:US:ICPSR or INFO:DDI:AGENCY:US:ICPSR. The latter can avoid  
174 confusion when these identifiers are not to be resolved via the DDI URN resolution  
175 mechanism.

### 176 **Compatibility with other URN resolution mechanisms**

- 177
- The proposed resolution mechanism is different from any other resolution mechanisms that  
178 are used within the academic / archival community. The main difference is that the  
179 resolution requires interpretation of specific elements of the URN and application and that it  
180 needs an algorithm to construct the proper URL. Currently, the proposed solution assumes  
181 that this will be taken care of by the agents. To be compatible with other mechanisms, this  
182 complexity should be taken care of by a reusable service that is available to other (non DDI-  
183 aware) applications.
- 184
- A problem can arise when DDI resolution will be integrated with existing resolver, because  
185 these might not be able or willing to handle the high load.

### 186 **Service records in DNS**

187 The Agency's name server maintains and advertises a list of machines and services capable of  
188 resolving their own URNs. See RFC 27825 for the explanation of the snippet below for  
189 example:

```
190     $ORIGIN example.com.  
191     @           SOA server.example.com. root.example.com. (  
192                1995032001 3600 3600 604800 86400 )  
193                NS  server.example.com.  
194                NS  ns1.ip-provider.net.  
195                NS  ns2.ip-provider.net.  
196     ; foobar - use old-slow-box or new-fast-box if either is  
197     ; available, make three quarters of the logins go to  
198     ; new-fast-box.  
199     _foobar._tcp  SRV 0 1 9 old-slow-box.example.com.  
200                SRV 0 3 9 new-fast-box.example.com.  
201     ; if neither old-slow-box or new-fast-box is up, switch to  
202     ; using the sysadmin's box and the server  
203                SRV 1 0 9 sysadmins-box.example.com.  
204                SRV 1 0 9 server.example.com.  
205     server       A   172.30.79.10  
206     old-slow-box  A   172.30.79.11  
207     sysadmins-box A   172.30.79.12  
208     new-fast-box  A   172.30.79.13  
209     ; NO other services are supported
```

---

<sup>5</sup> <http://www.ietf.org/rfc/rfc2782.txt>





## Data Documentation Initiative

210           \*.\_tcp               SRV  0 0 0 .  
211           \*.\_udp               SRV  0 0 0 .  
212

### 213 **Hosting essential infrastructure**

214 All resolution depends on proper and quick functioning of the name server that hosts the agency  
215 identifiers. To ensure availability, DNS allows for replication of these records over the name  
216 servers that are hosted by participants of the DDI Alliance and/or other agencies.

217 If the number of agencies increases, an additional delegation via country codes can be  
218 considered.

### 219 **DDI request API**

220 Use of delegation implies scalability and manageability as well as high availability. The  
221 insertions of extra parameters in the DDI query API (like auto-incremented TTL, access KEY,  
222 etc.) will prevent infinite circular references as well as grant access rights for the case that not  
223 all metadata is accessible for everyone.

### 224 **Resolution by algorithm**

225 Many persistent identifier experts claim that identifiers should not contain any information in  
226 their string. Although maybe not ideal, in this case it is a practical solution. It only implies that  
227 both the algorithm and the information should at least be persistent. The integration of  
228 versioning should assure this.

### 229 **Nested DDI objects**

230 When nesting schemes, best practice is to use the object and its parent Maintainable object,  
231 which provides a unique ID. The URN for the resource must not change, even though the  
232 scheme might be re-used in a different context.

233  
234 urn:ddi\_3\_0:       DataCollection   .       Coding  
235 =       MPC:        DC\_1[3.0]       .       PE\_2[1.0]       .       Code\_5

236 The definition of DDI URN does not include all object types. The more general case arising from  
237 the example in line 623 of DDI 3.0 Part 2 User Guide may contain an equation as follows:

238  $X.Y = A.(B_1...B_n).C$

239 This case is resolved as follows. Introduce in the DDI query API a variable called “da” to stand  
240 for disambiguator. The query request will then look like this:

241  $?X=B\&da=B_1...B_n\&Z=C$

242 The name disambiguator should also now be clear.

### 243 **Transforming a URN into a query string**



## Data Documentation Initiative

244 The process describes how the DDI-application transforms a URN into a query string. The DDI  
245 provider can also do this transformation. This would simplify the API and move the 'complexity'  
246 to the more centralized DDI provider. The drawback is that this doesn't allow querying for, e.g.,  
247 all variables within a variable scheme. The exact API will have to be determined with  
248 developers of DDI applications and DDI providers.

### 249 **3 References**

250 See document for embedded references.

#### 251 **3.1 Normative**

252 [RFC2119] S. Bradner, Key words for use in RFCs to Indicate Requirement Levels,  
253 <http://www.ietf.org/rfc/rfc2119.txt>, IETF RFC 2119, March 1997.

254 OASIS, Best Practice, [http://www.oasis-open.org/committees/uddi-](http://www.oasis-open.org/committees/uddi-spec/doc/bp/uddi-spec-tc-bp-template.doc)  
255 [spec/doc/bp/uddi-spec-tc-bp-template.doc](http://www.oasis-open.org/committees/uddi-spec/doc/bp/uddi-spec-tc-bp-template.doc), 2003



Data Documentation Initiative

256 **Appendix A. Acknowledgments**

257 The following individuals were members of the DDI Expert Workshop held 10-14 November  
258 2008 at Schloss Dagstuhl, Leibniz Center for Informatics, in Wadern, Germany.

259 Nikos Askitas, Institute for the Study of Labor (IZA)

260 Karl Dinkelmann, University of Michigan

261 Michelle Edwards, University of Guelph

262 Janet Eisenhauer, University of Wisconsin

263 Jane Fry, Carleton University

264 Peter Granda, Inter-university Consortium for Political and Social Research (ICPSR)

265 Arofan Gregory, Open Data Foundation

266 Rob Grim, Tilburg University

267 Pascal Heus, Open Data Foundation

268 Maarten Hoogerwerf, Data Archiving and Networked Services (DANS)

269 Chuck Humphrey, University of Alberta

270 Jeremy Iverson, Algenta Technology

271 Jannik Vestergaard Jensen, Danish Data Archive (DDA)

272 Kirstine Kolsrud, Norwegian Social Science Data Services (NSD)

273 Stefan Kramer, Yale University

274 Jenny Linnerud, Statistics Norway

275 Hans Jørgen Marker, Danish Data Archive (DDA)

276 Ken Miller, United Kingdom Data Archive (UKDA)

277 Meinhard Moschner, GESIS - Leibniz Institute for the Social Sciences

278 Ron Nakao, Stanford University

279 Sigbjørn Revheim, Norwegian Social Science Data Services (NSD)

280 Wendy Thomas, University of Minnesota

281 Mary Vardigan, Inter-university Consortium for Political and Social Research (ICPSR)



Data Documentation Initiative

- 282 Joachim Wackerow, GESIS - Leibniz Institute for the Social Sciences
- 283 Wolfgang Zenk-Möltgen, GESIS - Leibniz Institute for the Social Sciences
- 284



285 **Appendix B. Revision History**  
286

Rev	Date	By Whom	What
0.9	2009-03-21	Stefan Kramer	Began revision history tracking.



287

## 288 **Appendix C. Legal Notices**

289 Copyright © DDI Alliance 2009, *All Rights Reserved*

290

291 <http://www.ddialliance.org/>

292

293 Content of this document is licensed under a Creative Commons License:

294 Attribution-Noncommercial-Share Alike 3.0 United States

295

296 This is a human-readable summary of the Legal Code (the full license).

297 <http://creativecommons.org/licenses/by-nc-sa/3.0/us/>

298

299 You are free:

300 

- to Share - to copy, distribute, display, and perform the work

301 

- to Remix - to make derivative works

302

303 Under the following conditions:

304 

- Attribution. You must attribute the work in the manner specified by the author or

305 licensor (but not in any way that suggests that they endorse you or your use of the

306 work).

307 

- Noncommercial. You may not use this work for commercial purposes.

308 

- Share Alike. If you alter, transform, or build upon this work, you may distribute the

309 resulting work only under the same or similar license to this one. For any reuse or

310 distribution, you must make clear to others the license terms of this work. The best

311 way to do this is with a link to this Web page.

312 

- Any of the above conditions can be waived if you get permission from the copyright

313 holder.

314 

- Apart from the remix rights granted under this license, nothing in this license impairs

315 or restricts the author's moral rights.

316

317 **Disclaimer**

318

319 The Commons Deed is not a license. It is simply a handy reference for understanding the Legal Code

320 (the full license) — it is a human-readable expression of some of its key terms. Think of it as the user-

321 friendly interface to the Legal Code beneath. This Deed itself has no legal value, and its contents do not

322 appear in the actual license.

323

324 Creative Commons is not a law firm and does not provide legal services. Distributing of, displaying of, or

325 linking to this Commons Deed does not create an attorney-client relationship.

326 Your fair use and other rights are in no way affected by the above.

327

328 **Legal Code:**

329 <http://creativecommons.org/licenses/by-nc-sa/3.0/us/legalcode>

