

# Survey of EAI and IDN implementation

Irena Talašová\*



## Abstract

The main topic of this work is to find and summarize the current state of implementation and support of Unicode characters in the most known network protocols and applications. This is a topic that, in my opinion, deserves more attention from IT professionals and the general public, as it is likely to be the direction that application development will take. The following chapters summarize a brief introduction to this issue and will then introduce the current state of support for UTF-8 characters in selected network protocols and applications. Since, in general, the state of implementation and actual functionality of applications do not match the RFC specification of the corresponding network protocols, it is necessary to test and explore available applications and obtain relevant information. The information presented here are obtained both from expert articles and on-line accessible websites of organizations, and above all the results obtained during the testing. The information found can be used by all who want or plan to implement local language support in their developed tools. PCAP files containing accented network identifiers are provided as supplementary material that can be used for testing purposes.

**Keywords:** EAI — IDN — Unicode — Email address — Internet services — UA — Network protocols

**Supplementary Material:** [Demonstrative PCAPs](#)

\*xtalas04@stud.fit.vutbr.cz, Faculty of Information Technology, Brno University of Technology

## 1. Introduction

Increasing user's demand for multilingual Internet has led to the implementation of UTF-8 character support for some network protocols and applications. We are increasingly confronted with terms like EAI (Email Address Internationalization) in emails or IDN (Internationalized Domain Names) in the landscape of domain names. These principles also apply to other IT

areas and protocols such as SIP or FTP.

This creates a potential vulnerability for lawful interception and other software monitoring network communications, which can not process and recognize a network flow containing national characters. The aim of this work is to carry out research on the current state of implementation of IDN and EAI. It is not possible to examine all network applications, so the most

9

10  
11  
12  
13  
14  
15  
16

17 prominent representatives for whom the support will  
18 be tested are selected for each protocol being exam-  
19 ined. The knowledge gained is usable for all network  
20 application developers who are in the process of decid-  
21 ing whether to implement EAI and IDN support in the  
22 software.

23 The second chapter provides an introduction to  
24 the topic and explains the important concepts found  
25 in other parts of the paper. The third chapter focuses  
26 on IDN TLDs. The fourth chapter deals with practical  
27 knowledge with focus on SMTP, POP3 / IMAP, SIP  
28 and FTP protocols.

## 29 2. Preliminaries and Definitions

30 **Universal Acceptance (UA)** is the state where all do-  
31 main names and email addresses are accepted, vali-  
32 dated, stored, processed and displayed correctly and  
33 consistently by all applications, devices and systems [1],  
34 so internet applications and systems must treat all  
35 TLDs in a consistent manner, including new gTLDs  
36 and internationalized TLDs. This is a key measure  
37 of the success of IDNs (UA). Universal Acceptance  
38 is a foundational requirement for a truly multilingual  
39 Internet, in which users around the world can navigate  
40 entirely in local languages [2].

41 It is also an opportunity to support consumer choice  
42 and innovation in the domain name industry. Many  
43 systems do not recognize or appropriately process new  
44 domain names, primarily because the top-level domain  
45 may be new, more than three characters in length or in  
46 non-ASCII format (Internationalized Domain Names,  
47 or IDNs). The same is true for email addresses that  
48 incorporate these new domain names or use Unicode  
49 in the mailbox names [3].

50 **The Universal Acceptance Steering Group** [2],  
51 supported by Internet Corporation for Assigned Names  
52 and Numbers (ICANN), is a community-led initiative  
53 working on creating awareness and identifying and  
54 resolving problems associated with Universal Accep-  
55 tance. The purpose of these efforts is to help ensure  
56 a consistent and positive experience for Internet users  
57 globally. The group's primary objective is to help soft-  
58 ware developers and website owners understand how  
59 to update their systems to keep pace with an evolving  
60 domain name system (DNS).

61 Previous studies have shown that there are signif-  
62 icant barriers to Universal Acceptance of IDNs [4].  
63 Progress toward UA for IDNs is especially slow in ap-  
64 plications and security-related software. While there  
65 have been significant announcements of support for  
66 IDNs in email and other applications, the pace of up-  
67 take remains very slow [4].

68 **IDN** is an abbreviation for International Domain  
69 Names. This is a method and standard that allows  
70 people worldwide to use domain names in their local  
71 language. Domain name characters are encoded as  
72 UTF-8 string [5].

73 **EAI** is the extension that allows email addresses  
74 with IDNs in the domain part and/or Unicode (non-  
75 ASCII) characters in the mailbox name to function  
76 within the traditional email environment. Email soft-  
77 ware and services need to make specific changes to  
78 support EAI. Client software should display the do-  
79 main name and mailbox name in Unicode. Server  
80 software should confirm EAI-readiness (e.g. advertise  
81 SMTPUTF8 support) when making connection to an-  
82 other MTA [3]. An example of Internationalized Email  
83 Address is "háčky@čárky.cz".

84 **Punycode** is a way to represent IDNs with the lim-  
85 ited character set (A-Z, 0-9) supported by the domain  
86 name system [6]. For example, "háčkyčárky.cz" is  
87 encoded as "xn-hkyrky-ptac70bc.cz".

## 3. EAI & IDN in domain names

### 3.1 IDNs in the .cz TLD

88 CZ.NIC is the registry manager of the .cz country code  
89 top-level domains. According to their statement [7],  
90 they are technically prepared for the introduction of  
91 diacritics in Czech domains. The fact that the support  
92 has not yet been implemented is due to the alleged  
93 lack of interest among Czech users and organizations.  
94 On the other hand, second-level domains do not have  
95 problems with diacritics. Every two years the orga-  
96 nization conducts research of interest among Czech  
97 users. User's interest in IDNs is represented in Fig-  
98 ure 1. Numbers are given in percent.  
99  
100

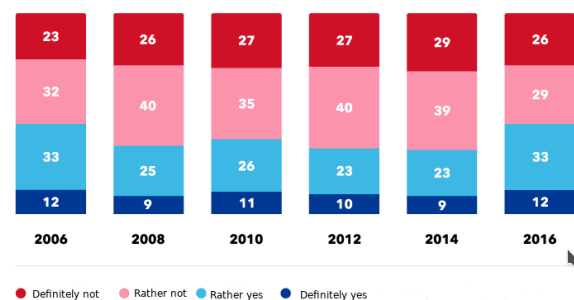


Figure 1. User's opinion on the introduction of diacritics in .cz domain [7].

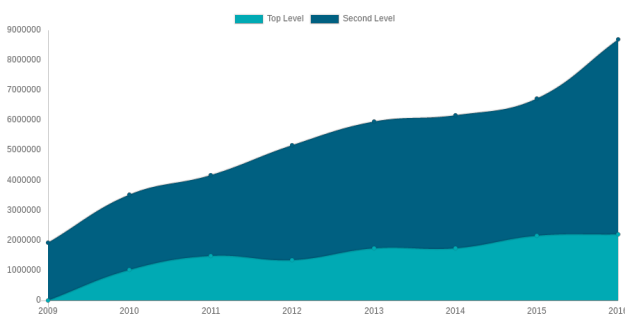
101 The introduction of IDN is a relatively complex  
102 process which, in addition to its benefits, includes sev-  
103 eral relatively important negative effects. According  
104 to CZ.NIC surveys, Czech users and organizations  
105 rank among the advantages of introducing IDN, clarity,  
106 simplicity, better readability and grammar and also it

107 can be used to protect their brand and it's quite com- 135  
 108 mon abroad. On the other hand, the disadvantages 136  
 109 are mainly problematic communication with abroad, it 137  
 110 can be confusing and there is possibility of errors. We 138  
 111 must consider some IDN security issues [7]. 139

### 112 3.2 IDNs in the other TLDs

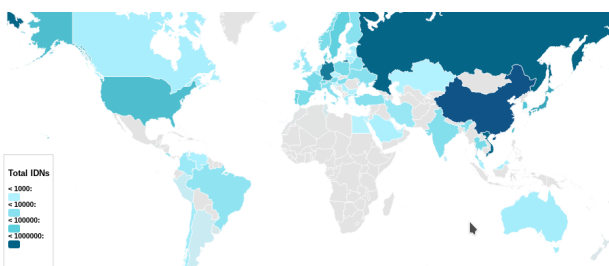
113 IDNs are only possible with some TLDs as .com, .eu, 142  
 114 .de, .pl, .at, .biz and others. One can learn more about 143  
 115 introducing national alphabet characters into .eu do- 144  
 116 mains at EURid website, a .eu domain administra- 145  
 117 tor [8]. 146

118 As can be seen in Figure 2 [9], since 2009, the 147  
 119 number of registered IDN domains has been growing 148  
 120 steadily. A particularly sharp increase is at second level 149  
 121 IDNs between 2015-2016. Nearly all of the growth 150  
 122 during that period is accounted for by .cn. 151



123 **Figure 2.** Growth of IDNs by level in the world [9].

124 World Map Growth of IDNs can be seen in Fig- 160  
 125 ure 3. This map shows the growth and distribution of 161  
 126 IDNs from 2013 onwards. The distribution is mapped 162  
 127 to the country of the host rather than the country of the 163  
 128 registrant. 164



129 **Figure 3.** Total number of registered IDNs in the world [9].

130 Other interesting charts can be found on [9].

### 129 3.3 IDN security issues

130 There are some security issues, the most common ones 174  
 131 are [7]:

- 132 1. **Typo-squatting** - the attacker registers the do- 175  
 133 main very similar to some of the most visited do- 176  
 134 mains. For example, attacker place pornography 177

138 or a fake bank page on his website, waiting for 139  
 139 users to get here by mistake and get some sen- 140  
 140 sitive data. To avoid the risk of typo-squatting, 141  
 141 you can register multiple domains. One can reg- 142  
 142 ister all combinations of common domains and 143  
 143 accented domains, but there are many. The inter- 144  
 144 national system for resolving disputes between 145  
 145 trade mark owners and domain name registrants 146  
 146 is described in [10]. 147

- 148 2. **DN homograph attack** is based on the inter- 149  
 149 changeability of some similar looking charac- 150  
 150 ters. This means that different interpretations of 151  
 151 these similar characters confuse the user. An ex- 152  
 152 ample is the "o" and "о" that looks the same, 153  
 153 although it can be written both in Latin and 154  
 154 Cyrillic. We get quite different from what we 155  
 155 wanted, though we clicked on a link that is seem- 156  
 156 ingly right. Homograph can also be related to 157  
 157 the Latin (e.g. "I" and "И", "O" and "О"). 158

### 3.4 Browsers

159 Browsers continue to be a bright spot for the use, dis- 160  
 160 play and resolution of IDNs [4]. Progress in browsers 161  
 161 has been made steadily in the last years. All major 162  
 162 web browsers support IDN. 163

## 4. EAI & IDN status survey

164 This chapter describes the differences between the 165  
 165 RFC protocol specifications and the tested behavior 166  
 166 of the applications. For the testing purposes, tools 167  
 167 listed below were installed in a Ubuntu-based virtual 168  
 168 machine. 169

### 4.1 Email protocols

170 The original **SMTP specification**, as with many other 171  
 171 protocols, did not support UTF-8 characters. The 172  
 172 change brought up SMTP extensions that allow UTF-8 173  
 173 characters in email headers and addresses. UTF-8 char- 174  
 174 acters are supported in both, the user part of the address 175  
 175 and the domain part. The domain part is converted to 176  
 176 a string without diacritics using the Punycode, when 177  
 177 formulating the DNS query. 178

179 At this time, none of the tested providers of **freemail** 180  
 180 **services**, such as Seznam, Atlas, Centrum, Volny, Epošta, 181  
 181 Yahoo, Gmail and Hotmail, has been able to create an 182  
 182 accented email box. The use of accented characters 183  
 183 occurring in the Czech alphabet such as "š", "č" or "ř" 184  
 184 leads to an error: "That email address contains in- 185  
 185 valid characters", when creating an e-mail box. Only 186  
 186 Gmail support both, IDN and EAI. It's interesting, that 187  
 187 mobile versions of Centrum.cz, Atlas.cz and Volny.cz 188  
 188 support IDN only. A message can be successfully sent 189

184 to "testmail@háčkyčárky.cz" and then the "IDN OK"  
185 returns. I think it is rather an implementation error.

186 As far as **SMTP servers** are concerned, the situa-  
187 tion is better. During testing, Postfix (3.1.0) and Exim  
188 (4.88) were able to work (send or receive) with EAI  
189 emails (EAI in header, such as "TO:háčky@čárky.cz").  
190 Only older servers Qmail (1.06) and Sendmail (8.15.2)  
191 did not support either IDN or EA.

192 Just like the original RFC for **SMTP, POP3/IMAP**  
193 RFC originally only counted with ASCII characters,  
194 EAI support was added later.

195 The first of the POP3 / IMAP server tests focused  
196 on finding and manipulating accented emails in email  
197 headers. Servers Dovecot (2.2.29), MS Exchange  
198 Server (2016) and Courier (0.78.2) passed this test  
199 seamlessly. Only Cyrus (2.4) could not find emails.

200 The second test was to create an accented account.  
201 I tried to create an account "Talašová@talasova.cz".  
202 Courier, Cyrus, and Dovecot were able to create this  
203 account. Only MS Exchange declined account creation  
204 due to unauthorized characters.

205 The third test was to log on to the server using the  
206 created account from the previous test. On Dovecot,  
207 MS Exchange Server and Courier servers, I was able  
208 to sign up for this account. Cyrus server only failed to  
209 sign in.

210 Among the best POP3 / IMAP clients in EAI sup-  
211 port are Outlook 2016 and Mutt, who fully support  
212 EAI. The remaining tested clients support only IDN.

213 I first tested whether clients can send an email to  
214 an address without diacritics from the address with  
215 diacritics (the diacritics are located in the FROM field  
216 only). Thunderbird (45.8.0), Outlook 2016 and Mutt  
217 (1.5.24) allowed to send such an email. In contrast, in  
218 the Roundcube Webmail (1.2.4) client, such message  
219 can not be sent.

220 In the second test, I focused on testing IDN sup-  
221 port. I was trying to send a message to autoresponder  
222 e-mail address: "testmail@háčkyčárky.cz". Thunder-  
223 bird, Roundcube Webmail, Pošta, Outlook, and Mutt  
224 successfully submitted e-mail.

225 The third and final test was focused on testing sup-  
226 port for EAI, that is, support for address with diacritics  
227 in the user section. I tried to send the message to  
228 "talašová@talasova.cz". Sending in Thunderbird and  
229 Roundcube Webmail was unsuccessful. Clients Mutt  
230 and Outlook 2016 allowed this message to be sent. The  
231 support status of UTF-8 strings is different for differ-  
232 ent clients. All tested clients supported the diacritics in  
233 the domain name. However, UTF-8 characters in the  
234 user part are mostly problem, especially in the "TO"  
235 field.

The tests **results** show the general effort of appli- 236  
cation developers to implement EAI support and be 237  
prepared if the demand for the diacritics increases. In 238  
general, EAI support in e-mails and web browsers is 239  
quite good. 240

## 4.2 File Transfer Protocol 241

**FTP** is one of the oldest and widely used Internet 242  
protocols. The original character set of this protocol 243  
was 7 bit ASCII [11]. Nowadays, however, there is 244  
a need to use national characters. The original FTP 245  
specification (RFC 2640) has been extended to support 246  
different character sets. Not ASCII mode but binary 247  
mode is used for transferring text files containing na- 248  
tional characters. UTF-8 characters within the file 249  
are not a problem for FTP. I tested Dolphin (4.14.3), 250  
FileZilla (3.26.2), FireFTP (2.0), gFTP (2.0.19), Kon- 251  
queror (4.14.3), Total Commander (9.0a), WS\_FTP 252  
(12.5), e-FTP (3.2.3.112) and Cerberus (8). 253

We will now focus on the use of diacritics in file 254  
names transmitted via FTP clients and servers. There 255  
is no problem with diacritics but with encoding. Trans- 256  
ferring the file with accented name was not a problem 257  
with any of the clients, if the file was created and 258  
subsequently moved with the help of the same client. 259  
Encoding errors sometimes occurred, when different 260  
clients were used. For example, using one client to 261  
save the file to the server and use the other one to view 262  
or download this file. Clients have support for different 263  
encoding of file names. The safest is to always use 264  
UTF-8 encoding. 265

## 4.3 Session Initiation Protocol 266

**SIP** is a text-based protocol that uses UTF-8 encoding. 267  
The field format in the header is defined by header 268  
name. It can be a sequence of UTF-8 octets or a com- 269  
bination of white characters, tokens, delimiters, and 270  
quoted strings. Message bodies may use different en- 271  
codings, this encoding should be reported in the value 272  
of the Content-Type header field. SIP messages may 273  
contain binary bodies; UTF-8 is assumed if the sender 274  
does not specify encoding. Caller names are encoded 275  
as UTF-8, so national characters can be used. UTF-8 276  
characters can only be used for descriptive character 277  
field values and those fields that are not expected to be 278  
interpreted by the parser [12]. 279

By RFC, the URI domain must contain only the 280  
selected ASCII characters. This is from the SIP URI 281  
syntax in the source [13]. However, for example, a 282  
Zoiper client sends a DNS query to the domain in the 283  
form given to it, whether punycode, hexadecimal, or 284  
UTF-8. Ekiga leaves the UTF-8 characters unchanged, 285

286 but the hexadecimal characters are changed to UTF-8,  
287 and punycode leaves also unchanged.

288 According to RFC, the user part of the SIP URI  
289 must contain only the selected ASCII characters or  
290 the hex hex sequence [13]. This is a system where  
291 "% xx" is written in place of the character, where x  
292 is the hexadecimal digits characterizing the bytes of  
293 the original character in UTF-8. In practice, each SIP  
294 client will handle differently an accented address. The  
295 Zoiper client accepts a user part in the form of talašová,  
296 tal% c5% a1ov% c3% a1 and tal% c5% a1. Accented  
297 characters are transmitted either in hexadecimal or as  
298 UTF-8 characters. While the Ekiga client change all  
299 these addresses to the hexadecimal form. In all cases,  
300 the Asterisk SIP server evaluated the address correctly.

#### 301 4.4 Other protocols

302 These improvements in e-mails, TLDs and browsers  
303 can not hide the fact that, in other parts of the Internet,  
304 Universal Acceptance is at best marginal and in some  
305 cases non-existent [4]. In the coming years we will see  
306 how and in what direction will EAI and UA support  
307 be developed within the Internet.

### 308 5. Conclusion

309 This work has highlighted modern tendencies towards  
310 the use of national languages when using network ser-  
311 vices, primarily web and email protocols. Not only  
312 was the goal to familiarize readers with concepts such  
313 as IDN, EAI or UA, but also to provide a comprehen-  
314 sive view of the results of testing and reviewing the  
315 current state of implementation of IDN and EAI sup-  
316 port in the Internet through various applications. The  
317 obtained information were used to develop module  
318 supporting national characters for lawful interception  
319 system. The main benefit is capture of data flows that  
320 include accents and security against attempts to hide  
321 communication using accented identifiers. And also  
322 general awareness raising about the possibilities, ben-  
323 efits and disadvantages of using UTF-8 characters in  
324 the Internet world.

### 325 Acknowledgements

326 I would like to thank my supervisor Ing. Libor Polčák,  
327 Ph.D. for his support.

### 328 References

329 [1] Don Hollander. Universal acceptance -  
330 an update, May 2016. [https://www.  
331 icann.org/en/system/files/files/  
332 gdd-summit-amsterdam-ua-17may16-en.  
333 pdf.](https://www.icann.org/en/system/files/files/gdd-summit-amsterdam-ua-17may16-en.pdf)

- [2] UASGtech. Universal acceptance, March 2018. 334  
<https://uasg.tech/home/>. 335
- [3] UASGtech. Quick guide to eai, March 2018. 336  
[https://uasg.tech/wp-content/  
337 uploads/2017/02/UASG014\\_  
338 20170206.pdf](https://uasg.tech/wp-content/uploads/2017/02/UASG014_20170206.pdf). 339
- [4] EURid. Universal acceptance, March 340  
2018. [https://idnworldreport.eu/  
341 year-2017/universal-acceptance/](https://idnworldreport.eu/year-2017/universal-acceptance/). 342
- [5] IETF. Internationalized domain names for appli- 343  
cations: Definitions and document framework, 344  
March 2018. [https://tools.ietf.org/  
345 html/rfc5890](https://tools.ietf.org/html/rfc5890). 346
- [6] Dynadot. What is punycode?, March 347  
2018. [https://www.dynadot.  
348 com/community/help/question/  
349 what-is-punycode](https://www.dynadot.com/community/help/question/what-is-punycode). 350
- [7] CZ.NIC. Cz.nic - idn - internationalized domain 351  
names, March 2018. [https://hackycarky.  
352 cz](https://hackycarky.cz). 353
- [8] EURid. Eurid, March 2018. [https://eurid.  
354 eu/cs/o-nas/eu-timeline/](https://eurid.eu/cs/o-nas/eu-timeline/). 355
- [9] EURid. Idn world report, March 2018. [https:  
356 //idnworldreport.eu](https://idnworldreport.eu). 357
- [10] David Lindsay. *International Domain Name Law: 358  
ICANN and the UDRP*. Hart Publishing, 2007. 359  
ISBN: 9781847313966. 360
- [11] IETF. Internationalization of the file transfer 361  
protocol, July 1999. [https://tools.ietf.  
362 org/html/rfc2640](https://tools.ietf.org/html/rfc2640). 363
- [12] Radhika Ranjan ROY. *Handbook on session ini- 364  
tiation protocol: networked multimedia commu- 365  
nications for IP telephony*. CRC Press, 2016. 366  
ISBN 9781498747707. 367
- [13] IETF. Sip: Session initiation protocol, March 368  
2018. [https://tools.ietf.org/html/  
369 rfc3261#section-25.1](https://tools.ietf.org/html/rfc3261#section-25.1). 370