# IRIs and IDNs: Problems of non-ASCII countries

Jakub Klímek

**American**
Standard
Code for
Information
Interchange

Since
1963

ASCII

IANA prefers the updated name
**US-ASCII**, which clarifies that this
system was developed in the **US**
and based on the typographical
symbols predominantly in use there

NUL SOH STX ETX EOT ENQ ▪ BEL BS HT LF VT FF CR SO SI ▪ DC1 DC2 DC3 DC4 NAK SYN EM ▪ ▪ ▪ ▪ ▪ FS GS RS US

!"#$%&'()*+,-./0123456789:;<=>?
@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]↑←

ACK ▪ ESC DEL

Meanwhile, in other countries...

ありがとうございます

Где ты сейчас?

Είμαι καλά

איך בין פֿײַן

죄송합니다

කොහොමද සැප සනීප?

ငါကစင်းပါပြ

# In Czechia, we are doing pretty good, but still...

Příliš žluťoučký kůň úpěl ďábelské ódy

# What did we do when we wanted to use the computers?

Příliš žluťoučký kůň úpěl ďábelské ódy

PÅ™Ã-liÅ¡ Å¾luÅ¥ouÄ□kÃ½ kÅ¯Å^ ÃºpÄ›l Ä□Ã¡belskÃ© Ã³dy

PrxEDlix
9A x9Elutouckx
FD kun xFApel dxE1belskxE9 xF3dy

Pý˘liç §luśouźkě k…Í ŁpŘl Ô belsk, ˘dy

| ISO 8859-2 | Windows CP1250 | Kamenický encoding | ... | ... | ... |

# What did we do when we wanted to use the computers?

Historically, all Czech geeks know that you **"just don't use"** diacritics in computers.

File names, Folders, SMS messages, URIs, ...

Příliš žluťoučký kůň úpěl ďábelské ódy

Prilis zlutoucky kun upel dabelske ody

# Unicode - UTF-8

Since
1993

Příliš žluťoučký kůň úpěl ďábelské ódy



Příliš žluťoučký kůň
úpěl ďábelské ódy.txt



15:22 ⚇ 🖂 🎤          ⏰ 🔇 📶 100% 🔋

‹   New conversation

Recipient                                    👤

+   Příliš žluťoučký kůň úpěl ďábelské ódy   ✈

|||          ◯          ‹

# Unicode identifiers: We really need IDNs and IRIs

1. Linked Data: Identifying all things with unique identifiers - IRIs
   a. Pieces of legislation -> stripping of diacritics leads to conflicts, loss of meaning, hard to read
   https://slovník.gov.cz/legislativní/sbírka/111/2009/pojem/má-úplné-označení-ustanovení-včetně-označení-právního-předpisu
   b. Cities (https://cs.wikipedia.org/wiki/Šumperk)
   c. Companies (https://čokoláda.cz / https://cokolada.cz)
   d. ...
2. Publishing data about things identified by IRIs
   a. Web pages
   b. XML Schemas
   c. JSON Schemas
   d. CSV on the Web schemas
   e. API endpoints
   f. ...

# IDN

https://háčkyčárky.cz/

Since 2003

RFC 5890

| Punycode ToASCII (RFC 3492) | → | ⬇ ⬆ | ← | Punycode ToUnicode (RFC 3492) |

https://xn--hkyrky-ptac70bc.cz/

# IRI

https://cs.wikipedia.org/wiki/Červená

Since 2005

RFC 3987

| Percent-encoding (RFC 3986) | → | ⬇ ⬆ | ← | Percent-decoding (RFC 3986) |

https://cs.wikipedia.org/wiki/%C4%8Cerven%C3%A1

URI → HTTP ← URI

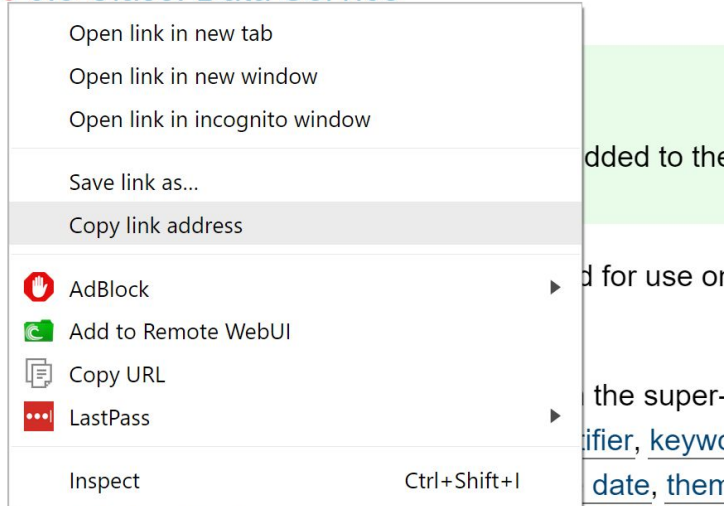# IDNs and IRIs: So far so good, right?

# How do you use URIs today?

# How do you use URIs today?

available.

See also: compression format.

§ 6.8 Class: Data Service

Open link in new tab
Open link in new window
Open link in incognito window

Save link as...
Copy link address

AdBlock
Add to Remote WebUI
Copy URL
LastPass

Inspect          Ctrl+Shift+I

qualified attribution

dded to the

d for use or

the super-
ifier, keyw
date, them

ctrl+v

Check this out:
https://www.instagram.com/

# OK, how about IDNs?

ctrl+c

CZ.NIC - IDN - International

🔒 xn--hkyrky-ptac70bc.cz

**CZ.NIC**

CZ.NIC - IDN - Internationalized ✕

🔒 https://háčkyčárky.cz

ctrl+v

Check this out: https://xn--hkyrky-ptac70bc.cz/

| Google Chrome | https://xn--hkyrky-ptac70bc.cz/ |
|---------------|----------------------------------|
| MS IE 11 | https://háčkyčárky.cz/ |
| MS Edge | https://háčkyčárky.cz/ |
| Firefox | https://háčkyčárky.cz/ |
| Vivaldi | https://xn--hkyrky-ptac70bc.cz |

# OK, how about IRIs?



| Google Chrome | https://cs.wikipedia.org/wiki/%C4%8Cerven%C3%A1 |
| MS IE 11 | https://cs.wikipedia.org/wiki/Červená |
| MS Edge | https://cs.wikipedia.org/wiki/Červená |
| Firefox | https://cs.wikipedia.org/wiki/%C4%8Cerven%C3%A1 |
| Vivaldi | https://cs.wikipedia.org/wiki/Červená |

# IDNs and IRIs combined?

ctrl+c
ctrl+v

https://háčkyčárky.cz/Červená

| Browser | IDN | IRI |
| --- | --- | --- |
| Google Chrome | xn--hkyrky-ptac70bc.cz | %C4%8Cerven%C3%A1 |
| MS IE 11 | háčkyčárky.cz | Červená |
| MS Edge | háčkyčárky.cz | Červená |
| Firefox | háčkyčárky.cz | %C4%8Cerven%C3%A1 |
| Vivaldi | xn--hkyrky-ptac70bc | Červená |

# Support of IRIs elsewhere: XML 1.0

```xml
<?xml version="1.0" encoding="UTF-8" ?>
<adresa
    xmlns="https://data.gov.cz/otevřené-formální-normy/adresy/"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation="https://data.gov.cz/otevřené-formální-normy/adresy/
    https://data.gov.cz/otevřené-formální-normy/adresy/draft/schémata/adresa.xsd">
```

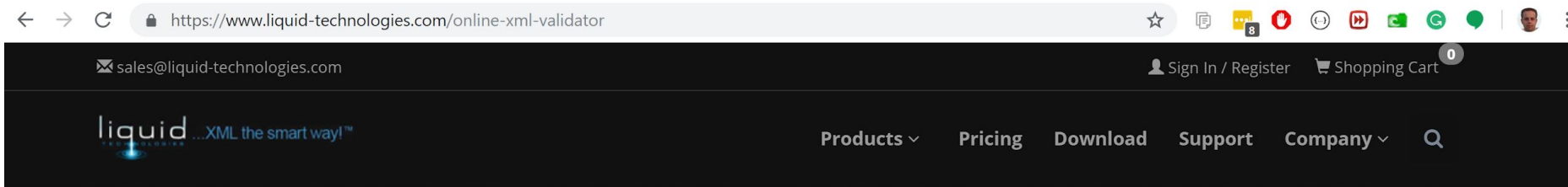https://data.gov.cz/otevřené-formální-normy/adresy/draft/příklady/0.xml

**This page contains the following errors:**

error on line 3 at column 64: xmlns: 'https://data.gov.cz/otevÅ□enÃ©-formÃ¡lnÃ-normy/adresy/' is not a valid URI

**Below is a rendering of the page up to the first error.**

# Support of IRIs elsewhere: XML 1.1

# Support of IRIs elsewhere: JSON Schema

```
{
    "$schema": "http://json-schema.org/draft-07/schema#",
    "$id": "https://data.gov.cz/otevřené-formální-normy/faktury/draft/schémata/faktury.json",
    "type": "array",
    "title": "Faktury",
    "items": {
        "$id": "#/items",
        "type": "object",
        "title": "Faktura",
        "properties": {
            ▶ "id": { … }, // 4 items
            ▶ "typ_dokladu": { … }, // 4 items
            ▼ "částka_bez_dph": {
                "$ref": "https://data.gov.cz/otevřené-formální-normy/základní-datové-typy/draft/schémata/částka.json"
            },
            ▼ "částka_s_dph": {
                "$ref": "https://data.gov.cz/otevřené-formální-normy/základní-datové-typy/draft/schémata/částka.json"
            },
            ▼ "částka_uhrazená": {
                "$ref": "https://data.gov.cz/otevřené-formální-normy/základní-datové-typy/draft/schémata/částka.json"
            },
```

Raw | Parsed

# Support of IRIs elsewhere: JSON Schema

# IRI support (i.e. Unicode in URIs) #59

⊘ **Open**   **awwright** opened this issue on Sep 18, 2016 · 3 comments

**awwright** commented on Sep 18, 2016    Member   ···

Most modern Web/hypermedia formats support IRIs instead of just URIs that are 7bit ASCII. IRIs are a superset of URIs that support full Unicode. For standards that only support URIs, IRIs have to be converted/escaped into a URI-compatible format.

👍 3

🏷 🔳 **awwright** added Type: Enhancement   Feedback period   labels on Sep 18, 2016

🔖 👤 **handrews** referenced this issue on Sep 19, 2016

### annotation: Multilingual meta data #53    ⊘ Closed

**handrews** commented on Sep 23, 2016    Member   ···

+1 for IRI support. It's 2016, let's be inclusive :-)

👍 2

---

## Assignees

No one assigned

## Labels

Type: Enhancement

core

## Projects

None yet

## Milestone

draft-future

2 participants

# And what about curl and IDN? (Ubuntu Linux)

```
klimek@KLIMEK-MFF-NTB:~$ curl -v -I https://háčkyčárky.cz/
*   Trying 217.31.205.51...
* Connected to háčkyčárky.cz (217.31.205.51) port 443 (#0)
*   subject: CN=xn--hkyrky-ptac70bc.cz
*   subjectAltName: host "háčkyčárky.cz" matched cert's "xn--hkyrky-ptac70bc.cz"

> HEAD / HTTP/2
> Host: xn--hkyrky-ptac70bc.cz
> User-Agent: curl/7.64.0
> Accept: */*
>

< HTTP/2 200
```

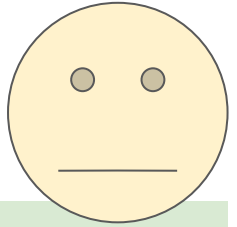# And what about curl and IDN? (Windows)

```
C:\Tools\curl-7.64.1-win64-mingw\bin>curl -v -I https://háčkyčárky.cz/

* Failed to convert háckycárky.cz to ACE;

* Closing connection -1

curl: (3) Failed to convert há
```

# And what about curl and IRI? (Ubuntu Linux)

```
klimek@KLIMEK-MFF-NTB:~$ curl -v -I --http1.1 https://cs.wikipedia.org/wiki/Červená
> HEAD /wiki/Červená HTTP/1.1
> Host: cs.wikipedia.org
> User-Agent: curl/7.64.0
> Accept: */*
>
< HTTP/1.1 200 OK
```

```
klimek@KLIMEK-MFF-NTB:~$ curl -v -I --http1.1
https://opendata-mvcr.github.io/podmínky-užití/není-chráněna-zvláštním-právem-pořizovatele-databáze/
> HEAD /podmínky-užití/není-chráněna-zvláštním-právem-pořizovatele-databáze/ HTTP/1.1
> Host: opendata-mvcr.github.io
> User-Agent: curl/7.64.0
> Accept: */*
>
< HTTP/1.1 400 Bad Request
```

# And what about curl and IRI? (Windows)

```
C:\Tools\curl-7.64.1-win64-mingw\bin>curl -v -I --http1.1 https://cs.wikipedia.org/wiki/Červená
> HEAD /wiki/Cervená HTTP/1.1
> Host: cs.wikipedia.org
> User-Agent: curl/7.64.1
> Accept: */*
>
< HTTP/1.1 404 Not Found
```

```
C:\Tools\curl-7.64.1-win64-mingw\bin>curl -v -I --http1.1
https://opendata-mvcr.github.io/podmínky-užití/není-chráněna-zvláštním-právem-pořizovatele-databáze/
> HEAD /podmínky-uzití/není-chránena-zvlástním-právem-porizovatele-databáze/ HTTP/1.1
> Host: opendata-mvcr.github.io
> User-Agent: curl/7.64.1
> Accept: */*
>
< HTTP/1.1 400 Bad Request
```

# Every time something goes wrong with IRIs:

I told you so!

Everyone knows that you **"just don't use"** diacritics in computers.

This should not be the message people get from using IRIs

# These problems are not new

MY URL ISN'T YOUR URL
- DANIEL STENBERG

https://daniel.haxx.se/blog/2016/05/11/my-url-isnt-your-url/

IMHO: IRIs, when used outside of HTTP (in clipboard, when shared, …), should be HUMAN readable. Agents are responsible for implementing the necessary conversions (IDNs & IRIs, not http://////////////////////////////example.com)