

# E-mail Address Reliability

**Maintenir la qualité des adresses e-mail**

**Vandy Berten**  
**Isabelle Boydens**  
Décembre 2013



# Management Summary (Français)



# Management Summary (Nederlands)



# Table des matières

<b>Management Summary (Français)</b>	<b>3</b>
<b>Management Summary (Nederlands)</b>	<b>5</b>
<b>Table des matières</b>	<b>7</b>
<b>1. Introduction</b>	<b>9</b>
<b>2. Aspects syntaxiques et techniques</b>	<b>12</b>
2.1. Décomposition d'une adresse	13
2.2. Vérification syntaxique	13
2.2.1. Syntaxe du nom de domaine	14
2.2.2. Syntaxe du nom d'utilisateur	14
2.2.3. Longueur maximale d'une adresse e-mail	15
2.3. Validation du « Top Level Domain »	17
2.4. Validation du nom de domaine	18
2.5. Validation d'une adresse	19
2.5.1. Serveur MX et protocole SMTP	19
2.5.2. Difficultés	21
2.5.3. Outils	22
2.6. Contrôle de consultation	23
2.6.1. Redirection de liens	24
2.6.2. Image avec identifiant unique	26
2.6.3. Contrôle lors d'autres contacts	27
2.6.4. Réponse aux e-mails	27
2.7. Indicateurs de qualité et statistiques	28
2.7.1. Erreurs syntaxiques	28
2.7.2. Dégressivité de la validité dans le temps	29
2.7.3. Dépendance au nom de domaine	30
2.7.4. Proportion professionnels-particuliers	31
<b>3. Stratégie de bonne gestion des adresses e-mail</b>	<b>34</b>
3.1. Arbitrage stratégique	34
3.1.1. Rejeter les mauvais ou accepter les bons ?	34
3.1.2. Accepter l'exotisme ?	35
3.1.3. Quel contrôle sur les serveurs d'envoi d'e-mail ?	35
3.2. Aspects syntaxiques	36
3.2.1. Catégorisation	36
3.2.2. Syntaxe spécifique	38
3.2.3. Suggestions de correction	38
3.3. Validation d'adresse	39
3.4. Suspicion d'erreurs	40

3.4.1. Par matching interne	40
3.4.2. Noms de domaine fréquents	42
3.5. Matching et dédoublement	42
3.6. Batch ou on-line ?	43
3.7. Historique de la validité dans le temps et monitoring	44
3.8. Traitement on-line	45
3.8.1. Mise en place de tests en entrée	45
3.8.2. Suivi de la validité dans le temps	46
3.9. Traitement batch des fichiers existants	48
3.10. Organisation	48
<b>4. Panorama d'outils existants sur le marché</b>	<b>50</b>
4.1. Vérificateurs syntaxiques	50
4.2. Testeurs d'existence	51
4.2.1. ServiceObjects	52
4.2.2. EmailVerify for .NET	52
4.3. Outils de suivi	53
4.4. Data Quality Tools	54
4.4.1. OpenRefine (Google refine)	54
4.4.2. IntoDQ (Trillium)	54
4.4.3. RedPoint	54
4.4.4. Human Inference	55
4.5. Outils CRM	55
<b>5. Conclusion</b>	<b>56</b>
<b>6. Bibliographie</b>	<b>58</b>
<b>7. Annexes</b>	<b>59</b>
7.1. Vérification syntaxique	59
7.1.1. Vérification syntaxique générale	59
7.1.2. Vérification syntaxique spécifique	59
7.2. Typologie des événements	61
7.3. Adresses e-mail officielles pour les citoyens	63
7.4. Éviter les risques liés au spam	65
<b>8. Glossaire</b>	<b>67</b>

# 1. Introduction

*egovernment*

*benchmarks*



*egovernment*

*benchmarks*

*bounce mails*

A

\*

A

\*

\*

A

✱

✱

✱

A



## 2. Aspects syntaxiques et techniques

---

## 2.1. Décomposition d'une adresse



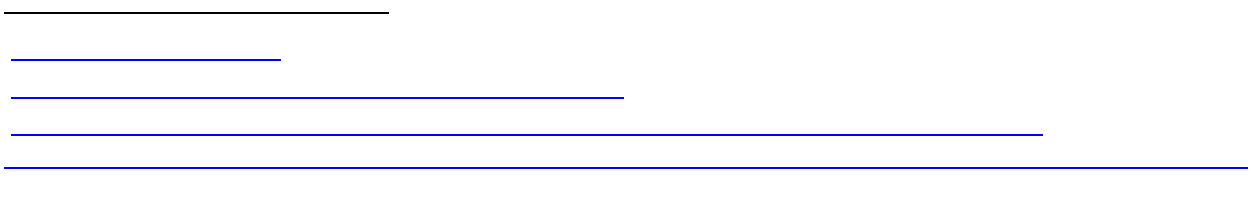
---

## 2.2. Vérification syntaxique

### 2.2.1. Syntaxe du nom de domaine



### 2.2.2. Syntaxe du nom d'utilisateur





*underscore)*

### 2.2.3. Longueur maximale d'une adresse e-mail

- 

- 

- 

- 

- 

- 

---

<http://tools.ietf.org/html/rfc821>

<http://tools.ietf.org/html/rfc5321>

<http://tools.ietf.org/html/rfc3696>

[http://www.rfc-editor.org/errata\\_search.php?rfc=3696&eid=1003](http://www.rfc-editor.org/errata_search.php?rfc=3696&eid=1003)

[http://www.rfc-editor.org/errata\\_search.php?rfc=3696&eid=1690](http://www.rfc-editor.org/errata_search.php?rfc=3696&eid=1690)

- 

abcdefghijklmnopqrstuvwxyzabcdefghijklmnopqrstuvwxyzab  
cdefghijkl@abcdefghijklmnopqrstuvwxyz.abcdefghijklmnopqrstuvwxyz  
qrstuvwxyz.abcdefghijklmnopqrstuvwxyz.abcdefghijklmnopqrstuvwxyz  
qrstuvwxyz.abcdefghijklmnopqrstuvwxyz.abcdefghijklmnopqrstuvwxyz  
qrstuvwxyz.abcdefghijklmnopqrstuvwxyz.be

- 
- 
- 

---

## 2.3. Validation du « Top Level Domain »

14

15



சிங்கப்பூர்

中國  
الجزائر

---

## 2.4. Validation du nom de domaine



17

18

---

## 2.5. Validation d'une adresse

### 2.5.1. Serveur MX et protocole SMTP

```
C:\>nslookup -q=mx gmail.com
[...]
Non-authoritative answer:
gmail.com mail exchanger = 5 gmail-smtp-in.1.google.com.
gmail.com mail exchanger = 10 alt1.gmail-smtp-in.1.google.com.
gmail.com mail exchanger = 20 alt2.gmail-smtp-in.1.google.com.
[...]
```

```
C:\>telnet gmail-smtp-in.l.google.com. 25
Trying 173.194.78.26...
Connected to gmail-smtp-in.l.google.com.
[...]
EHLO bxl.mapetitesociete.be
250-mx.google.com at your service, [91.183.59.xxx]
[...]
MAIL FROM:<albert.leroy@bxl.mapetitesociete.be>
250 2.1.0 OK pn9si600796wjc.42 - gsmtip
RCPT TO:<leroy.mariecelestine@gmail.com>
550-5.1.1 The email account that you tried to reach does not exist.
[...]
RCPT TO:<mariecelestine.leroy@gmail.com>
250 2.1.5 OK pn9si600796wjc.42 - gsmtip
QUIT
221 2.0.0 closing connection pn9si600796wjc.42 - gsmtip
```

*bounce mail*

*hards*

*softs*

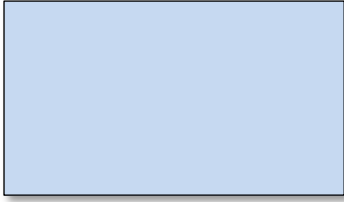
## 2.5.2. Difficultés



▪

▪

- 



*catch-all*

- 

*blacklister*

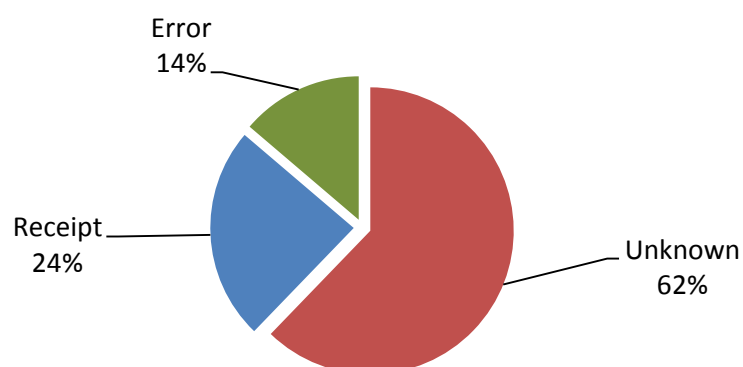


### 2.5.3. Outils



## 2.6. Contrôle de consultation

*: Exemple typique de situation suite à un envoi de campagne de communication : 24 % ont généré un accusé de réception, 14 % une erreur, et dans 62 % des cas, aucune information n'est disponible*



### 2.6.1. Redirection de liens

`http://mysite.com/track?m=albert@gmail.com&dst=www.smals.be`

`« albert@gmail.com;www.smals.be/a_page »`

`« YWxiZXJ0QGdtYWkuY29tO3d3dy5zbWFscy5iZS9hX3BhZ2U=»19`

`http://mysite.com/track?YWxi[...]t03d3dy5zbWFscy5iZS9hX3BhZ2U=`

```
<a href='http://mysite.com/track?YWxi[...]t03d3dy
5zbWFscy5iZS9hX3BhZ2U='>http://www.smals.be/a_page</a>
```

```
<html><head><meta http-equiv="refresh" content="0;
URL="http://www.smals.be/a_page"></head></html>
```

`www.smals.be/a_page`

▪

▪

▪

▪

▪



## **2.6.2. Image avec identifiant unique**

### **2.6.3. Contrôle lors d'autres contacts**

### **2.6.4. Réponse aux e-mails**

---

## 2.7. Indicateurs de qualité et statistiques

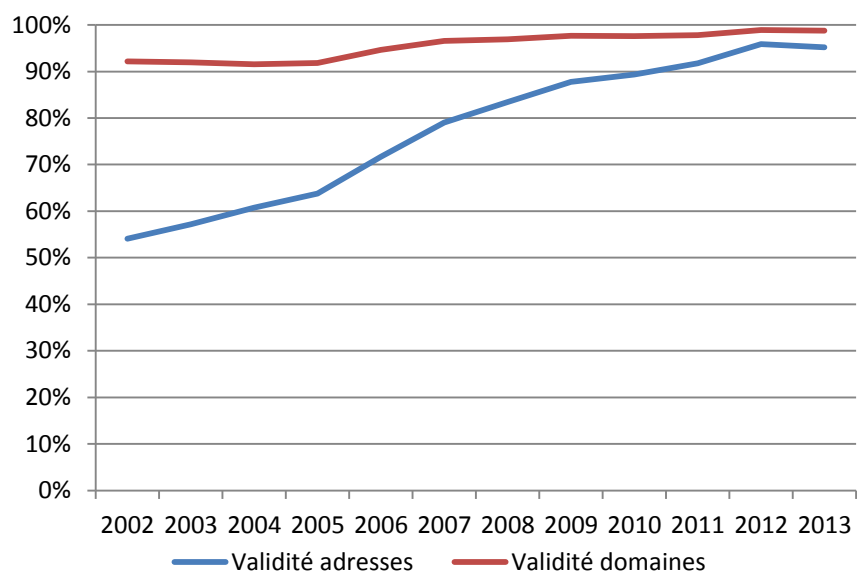
### 2.7.1. Erreurs syntaxiques

22

▪

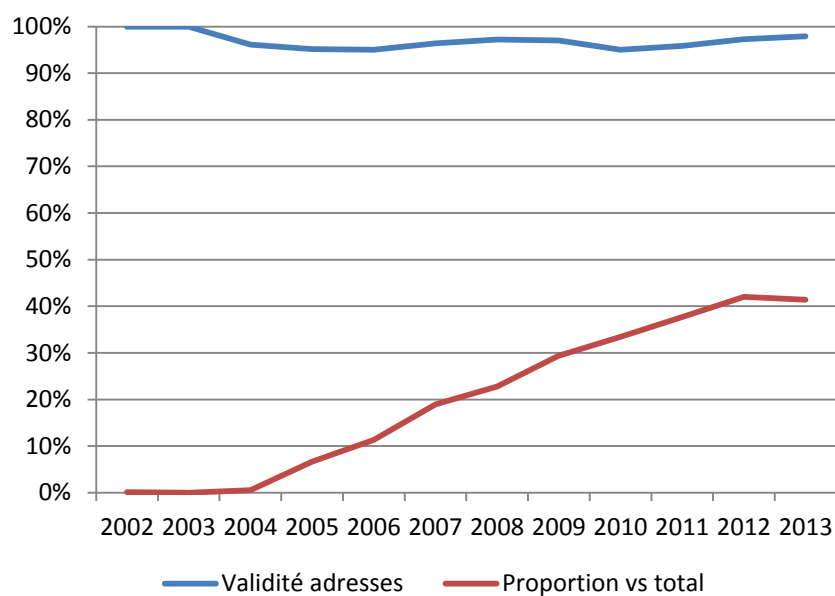
▪

*Validité des adresses e-mail et de leur nom de domaine associé d'une base de données de citoyens en fonction de l'année d'introduction*



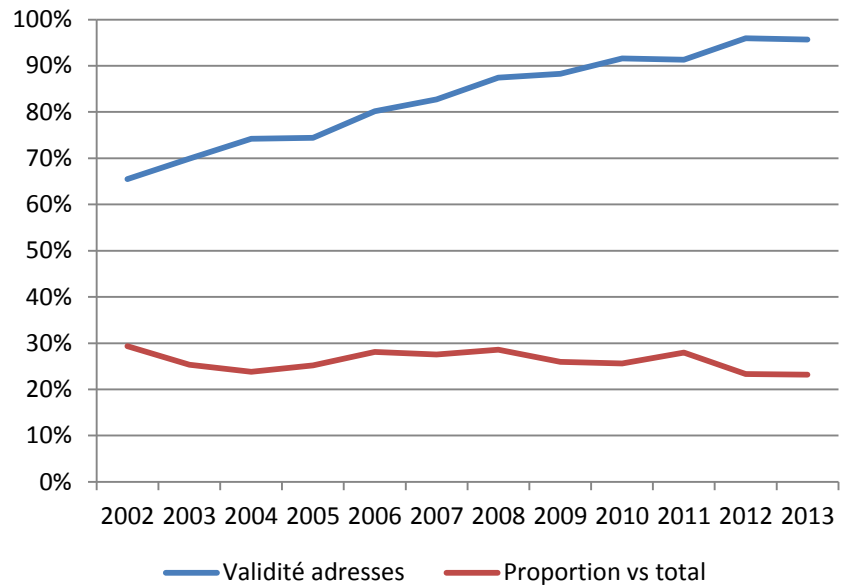
### 2.7.2. Dégressivité de la validité dans le temps

*Validité des  
adresses Gmail et  
proportion dans la base de  
données*



### 2.7.3. Dépendance au nom de domaine

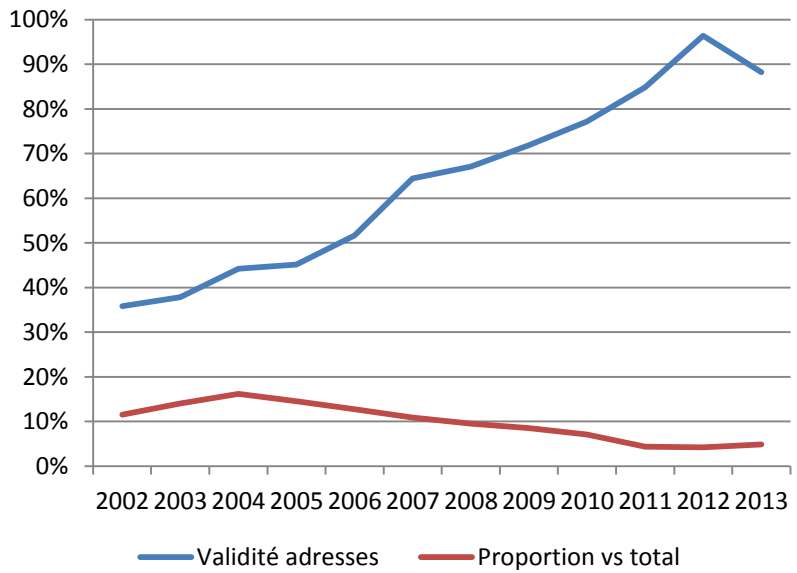
*Validité des  
adresses Hotmail et  
proportion dans la base de  
données*



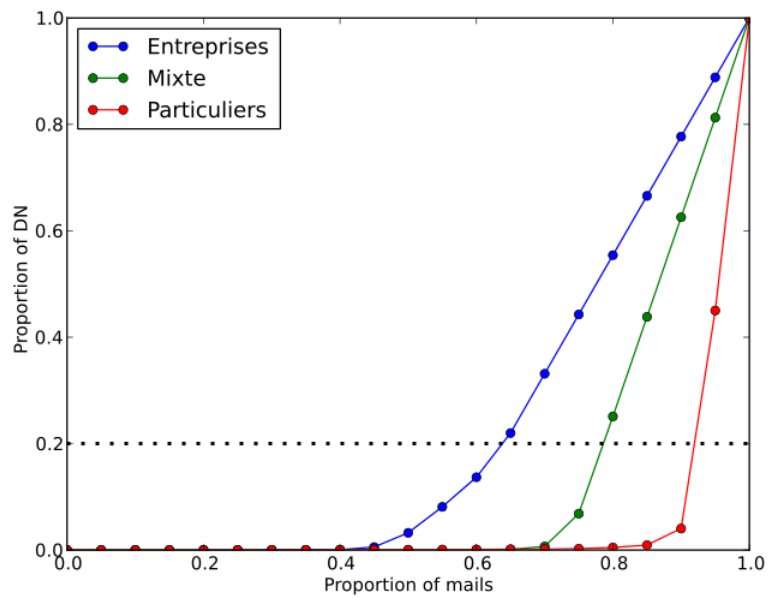
#### 2.7.4. Proportion professionnels-particuliers

- 
- 
-

*Validité des adresses  
Skynet et proportion dans la base  
de données*




*Distribution des  
adresses par rapport aux  
noms de domaine*



## 3. Stratégie de bonne gestion des adresses e-mail

---

### 3.1. Arbitrage stratégique

#### 3.1.1. Rejeter les mauvais ou accepter les bons ?

---

### 3.1.2. Accepter l'exotisme ?



- 
- 
- 
- 

### 3.1.3. Quel contrôle sur les serveurs d'envoi d'e-mail ?

---

---

---

## 3.2. Aspects syntaxiques

### 3.2.1. Catégorisation

- 
- 
- 

*Adresses certainement fausses*

`^[^@,;:\s]+@[^@,;:\s]+$`

$^{\wedge}[\wedge\text{@},;:\backslash\text{s}]+\text{@}([\wedge\text{@},;:\backslash\text{s}+.-]+([\wedge\text{@},;:\backslash\text{s}+.-]+)*)\text{\$}$

$[\wedge\text{@},;:\backslash\text{s}]+\text{@}([\backslash\text{p}\{\text{L}\}0-9]+([\wedge\text{@},;:\backslash\text{s}+.-]+)*)[\wedge\text{@},;:\backslash\text{s}+.-]+\text{\$}$

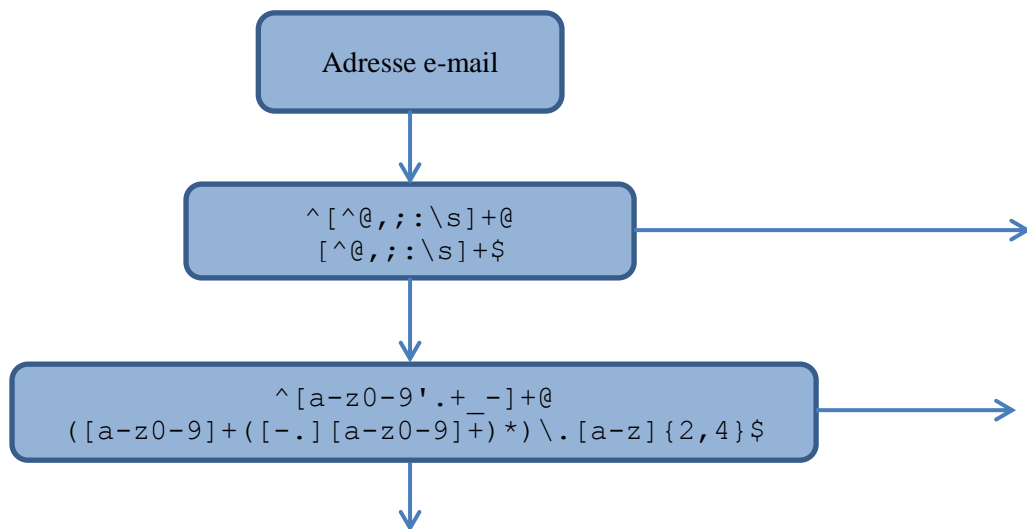
### **Adresses certainement correctes**

$^{\wedge}[\text{a-z}0-9'.+_-]+\text{@}([\text{a-z}0-9]+([\wedge\text{@},;:\backslash\text{s}+.-]+)*)\backslash\text{.}[\text{a-z}]{2,4}\text{\$}$

*underscore)*

$^{\wedge}[\text{a-z}0-9'.+_-]+([\wedge\text{@},;:\backslash\text{s}+.-]+([\wedge\text{@},;:\backslash\text{s}+.-]+)*)\backslash\text{.}[\text{a-z}]{2,4}\text{\$}$

## Adresses suspectives



### 3.2.2. Syntaxe spécifique

### 3.2.3. Suggestions de correction

- 
- - 
  -
- → →
- 
- 
- 

---

### 3.3. Validation d'adresse

---

## **3.4. Suspicion d'erreurs**

### **3.4.1. Par matching interne**

*underscore*

▪  
▪

### 3.4.2. Noms de domaine fréquents

			A	A

---

## 3.5. Matching et dédoublonnage

---

### **3.6. Batch ou on-line ?**

---

### **3.7. Historique de la validité dans le temps et monitoring**



---

## **3.8. Traitement on-line**

### **3.8.1. Mise en place de tests en entrée**

- 
- 
- 

### **3.8.2. Suivi de la validité dans le temps**

- 
- 
- 
- 

- 

- 

-

---

## 3.9. Traitement batch des fichiers existants

- 
  
- 
  
- - 
  - 
  - 
  
  - 
  -
  
- 
  
- 

---

## 3.10. Organisation

*bounces*

*return on investment*

- 
- 
- 
  
- 
-

## 4. Panorama d'outils existants sur le marché

---

### 4.1. Vérificateurs syntaxiques

- 

- 

- 

-

- 

- 

```
^[a-zA-Z0-9.!#$%&'*/=?^_`{|}~-]+@  
[a-zA-Z0-9-]+(\.[a-zA-Z0-9-]+)*$
```

- 

---

## 4.2. Testeurs d'existence

- 

- 

- 

- 

-

#### **4.2.1. ServiceObjects**

---

#### **4.2.2. EmailVerify for .NET**

---

---

### 4.3. Outils de suivi

- 

- 

-

---

## **4.4. Data Quality Tools**

### **4.4.1. OpenRefine (Google refine)**

---

### **4.4.2. IntoDQ (Trillium)**

---

### **4.4.3. RedPoint**

---

#### 4.4.4. Human Inference

---

---

---

## 4.5. Outils CRM

*bounces*

## 5. Conclusion

*return on investment*



## 6. Bibliographie

- [1] D. Clément, B. Laboisse, D. Duquennoy et A. Micheaux, «Non qualité de données & CRM : quel coût ?,» chez *QDC 2008*.
- [2] I. Boydens, «Strategic Issues Relating to Data Quality for E-government: Learning from an Approach Adopted in Belgium,» *Practical Studies in E-Government : Best Practices from Around the World*, pp. 113-130 (chapitre 7), 2011.
- [3] Y. Bontemps, I. Boydens et D. Van Dromme, «Data Quality : tools,» Bruxelles, 2007.
- [4] I. Boydens, *Informatique, normes et temps*, Bruxelles: Bruylant, 1999.
- [5] I. Boydens, A. Hulstaert et D. Van Dromme, «Gestion intégrée des anomalies,» 2011. [En ligne]. Available: <http://www.smalsresearch.be/publications/document?docid=62>.

## 7. Annexes

---

### 7.1. Vérification syntaxique

#### 7.1.1. Vérification syntaxique générale

```
^[^@,;:\s]+@[([^@,;:\s+\.-]+([-\.]^[^@,;:\s+\.-]+)*)$
```

```
^[^@,;:\s]+@  
([\p{L}0-9]+([-]\p{L}0-9+)*[.])*\p{L}{2,}$
```

```
^[a-z0-9\!\+\_]+\.[a-z0-9\!\+\_]+@  
([a-z0-9]+([-\.] [a-z0-9]+)*)\.[a-z]{2,4}$
```

#### 7.1.2. Vérification syntaxique spécifique



## 7.2. Typologie des évènements


---

### 7.3. Adresses e-mail officielles pour les citoyens

---

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

- 
- 
- 
-

---

## 7.4. Éviter les risques liés au spam

*spam*

*bots*

*honeypot*

*bots*

*bots*

*honeypot*

*spamtraps*

*honeypot*

*spamtraps »*

## 8. Glossaire

✱