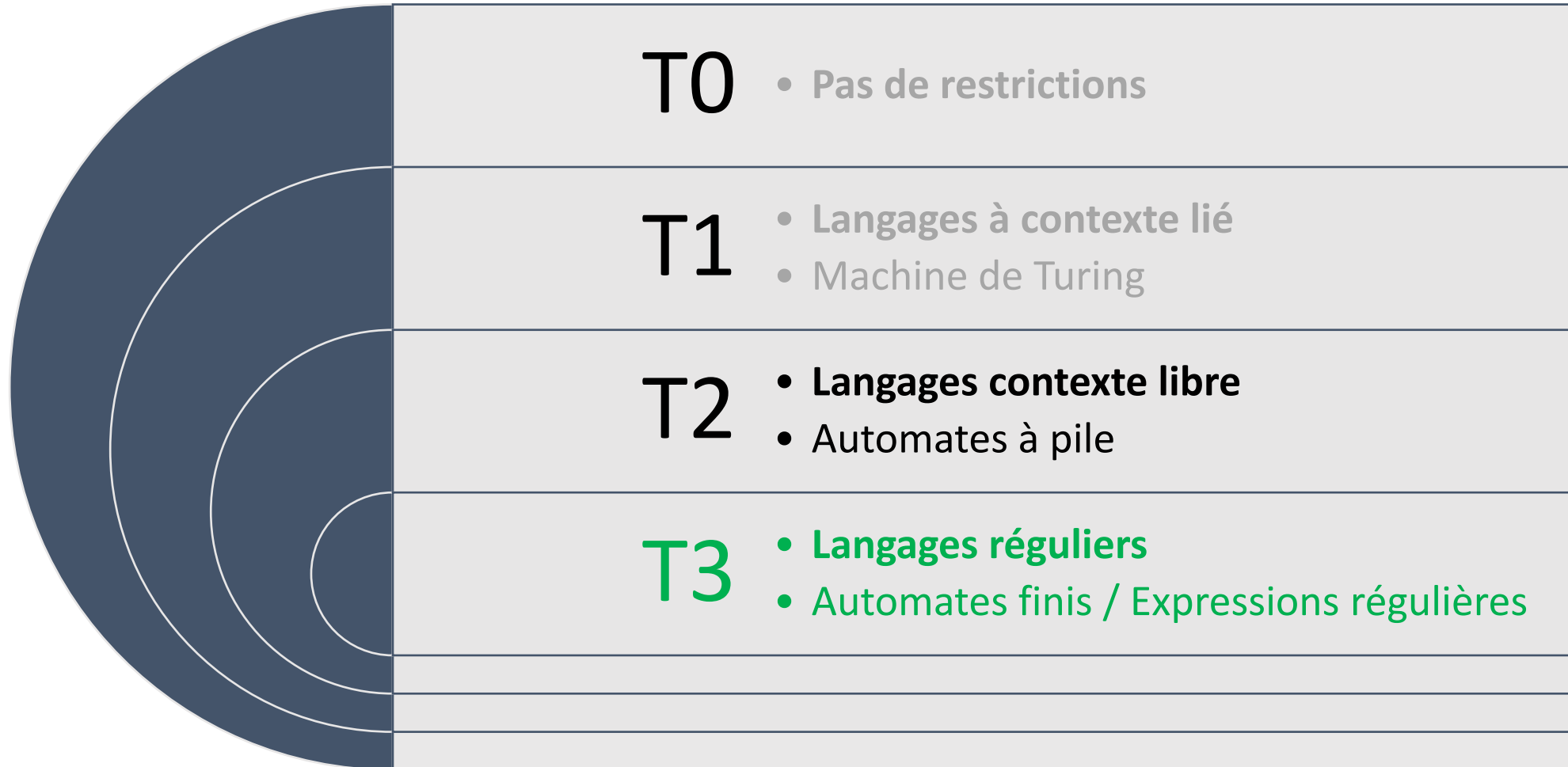




ISEN CIR ³

Théorie des langages

1. Expressions régulières



- **Grammaire de type 3** : Dans toutes les productions :
 - La partie gauche est composée d'un seul non-terminal
 - La partie droite est soit un terminal, soit une paire terminal/non-terminal,
 - en respectant la même position du non-terminal pour toutes les règles (à gauche ou à droite).
- Exemple d'une grammaire régulière avec les non-terminaux à droite:

$$S \rightarrow Aa \mid Sb \mid \varepsilon$$
$$A \rightarrow Sa \mid Ab$$

- Autre exemple, mais cette fois d'une grammaire régulière gauche:

$$S \rightarrow 0A \mid 1B \mid \varepsilon$$
$$A \rightarrow 0S \mid 1C$$
$$B \rightarrow 0C \mid 1S$$
$$C \rightarrow 0B \mid 1A$$

- Rappel: Cette grammaire permet de reconnaître des nombres binaires ayant un nombre paire de 0 et un nombre pair de 1

- Les langages réguliers sont la forme la plus simple des langages.
- Leur vérification par des outils informatiques est très simple et très rapide.

- Appelés aussi:

Expressions rationnelles

Motifs

Expressions normales

Regex

Regexp

Définition : Ensemble des langages réguliers

- L'ensemble des langages réguliers (ou langages rationnels), sur un alphabet A est le plus petit ensemble R des langages satisfaisant les conditions suivantes :

- $\emptyset \subset R$
- $\{\epsilon\} \subset R$
- $\forall a \in A : \{a\} \subset R$
- Si A et $B \subset R$ alors :

$$A.B \subset R$$

(La concaténation)

$$A \cup B \subset R$$

(Le choix)

$$A^* \subset R$$

(La répétition)

- Les Langages réguliers sont les langages définis en utilisant une association de ces trois opérateurs :

La concaténation

Le choix

La répétition

- Les langages réguliers peuvent être décrits de 3 façons:
 1. Une **grammaire régulière** (linéaire à droite ou à gauche).
 2. Une **expression régulière**.
 3. Un **automate fini** (Prochain cours).
- Pour démontrer qu'un langage est régulier, il suffit de le décrire par une des 3 façons.
- Il existe des méthodes (algorithmes) pour passer d'une forme à une autre.

- Une expression régulière est une formule qui combine les éléments de l'alphabet et un ensemble d'opérateurs :
 $\{ | \text{ (choix)}, \cdot \text{ (concaténation)}, * \text{ (répétition)}, (,), \epsilon \}.$
- Une expression régulière exprime un langage régulier d'une manière plus synthétique.
- Les opérations sur les expressions correspondent aux opérations sur les langages

(Choix)	\cup (union)
\cdot (Concaténation)	\cdot (Produit de langages)
* (Répétition 0 à n fois)	* (Etoile, Fermeture de Kleene)

- Constante naturelle en décimal

$0 \mid (1 \mid 2 \mid 3 \mid 4 \mid 6 \mid 7 \mid 8 \mid 9) \cdot (0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^*$

Simplifiée (dans certains langages) par : $0 \mid [1-9][0-9]^*$

- Binaires multiples de 2

$(0 \mid 1)^* \cdot 0$

- Identificateurs en C++

$(a \mid \dots \mid z \mid A \mid \dots \mid Z \mid _)\cdot(a \mid \dots \mid z \mid A \mid \dots \mid Z \mid _ \mid 0 \mid \dots \mid 9)^*$

- Langages de programmation

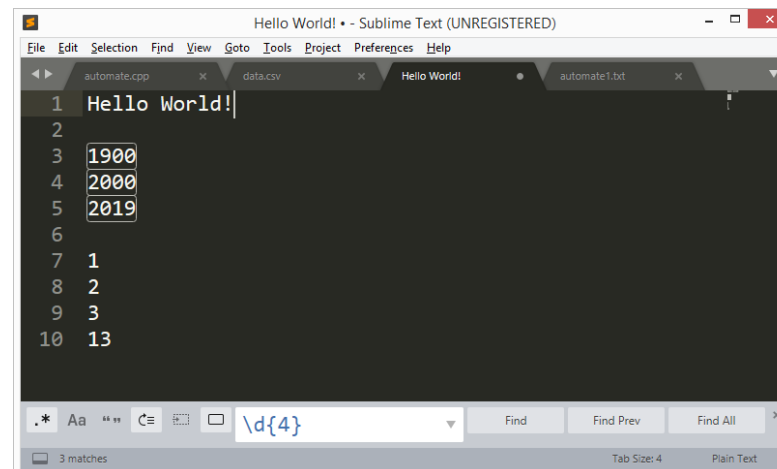
C++, Javascript, Perl, PHP, ...

- Editeurs de texte

Word, Notepad++, Sublime, ...

- Environnement de commande

grep

A screenshot of a question configuration interface. The 'Question Title' field contains 'Text Question'. The 'Question Type' is set to 'Text'. Under the 'Advanced settings' section, 'Data Validation' is checked. A 'Regular expression' is entered: /\d{2}[/-]\d{2}[/-]\d{4}/. The 'Required question' checkbox is unchecked.

```
alw@DESKTOP-FDCJ5J0:~$ egrep d\.g
He was a good dog
He was a good dog
alw@DESKTOP-FDCJ5J0:~$ egrep d\\.g
He was a good dog
This is the match string d.g
This is the match string d.g
alw@DESKTOP-FDCJ5J0:~$
```

- Il existe plusieurs variantes et de notations pour les expressions régulières.

- le standard C++2011 énumère les différentes notations:

Ecmascript, POSIX, awk, grep, ...

https://fr.wikipedia.org/wiki/Expression_r%C3%A9guli%C3%A8re

- Chacune de ses normes ajoute, en plus des trois opérateurs de base, des opérateurs relatifs aux classes de caractères (Chiffres, majuscules, minuscules, ...), des opérateurs relatifs aux répétitions et autres.

[]	Classe de caractères. Ex : [aeyuio] [0-9] [A-Z] [a-z]
.	Tout caractère sauf le caractère fin de ligne
*	Quantificateur : 0 ou n fois
+	Quantificateur : 1 ou n fois (au moins une fois)
?	Quantificateur : Au plus une occurrence (0 ou une fois)
{m, n}	Quantificateur : Entre n et m fois. Ex : {4, 8} : répétition entre 4 et 8 fois
^	Début de ligne
\$	Fin de ligne

REGEX SYNTAX	MEANING	EXAMPLE	MATCHES	DOES NOT MATCH
.	Any single character	go.gle	google, goggle	gogle
[abc]	Any of these character	analy[zs]e	analyse, analyze	analyxe
[a-z]	Any character in this range	demo[2-4]	demo2, demo3	demo1, demo5
[^abc]	None of these characters	analy[^zs]e	analyxe	analyse, analyze
[^a-z]	Not a character in this range	demo[^2-4]	demo1, demo5	demo2, demo3
	Or	demo example	demo, demos, example	test
^	Starts with	^demo	demos, demonstration	my demo
\$	Ends with	demo\$	my demo	demonstration
?	Zero or one times (greedy)	demos?123	demo123, demos123	demoA123
??	Zero or one times (lazy)			
*	Zero or more times (greedy)	goo*gle	gogle, goooogle	goggle
*?	Zero or more times (lazy)			
+	One or more times (greedy)	goo+gle	google, goooogle	gogle, goggle
+?	One or more times (lazy)			
{n}	n times exactly	w{3}	www	w, ww
{n,m}	from n to m times	a{4, 7}	aaaa, aaaaa, aaaaaa, aaaaaaa	aaaaaaaa, aaa, a
{n,}	at least n times	go{2,}gle	google, gooogle, goooogle	ggle, gogle
()	Group	^(demo example)[0-9]+	demo1, example4	demoexample2
(?:)	Passive group (Useful for filters)			
\	Escape	AU\\$10	AU\$10, AU\$100	AU10, 10
\s	White space			
\S	Non-white space			
\d	Digit character			
\D	Non-digit character			
\w	Word			
\W	Non-word (e.g. punctuation, spaces)			

- Username

`[a-z0-9_-]{3,16}`

- Password

`[a-zA-Z0-9_-]{6,18}`

- Couleurs hexadécimales en HTML/CSS

`#?([a-f0-9]{6}|[a-f0-9]{3})`

- Adresse e-mail

`([a-z0-9_\.-]+)@([\da-z_\.-]+\.[a-z\.]{2,6})`

- URL

`(https?:\//)?([\da-z\.-]+)\.([a-z\.]{2,6})([\w \.-]*)*\/?`

- Adresse IP

`((25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?)\.){3}(25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?)`

- Tag HTML

`<([a-z]+)([^\<]+)*(?:> (.*) <\/\1>|\s*\/>)`

<https://regex101.com/>

<http://regexpr.com/>

<https://www.regexpal.com/>