

# Enumeración de anotaciones con delay constante

Clase 27

IIC 2223

Prof. Cristian Riveros

# Autómata con anotaciones (recordatorio)

## Definición

Un autómata con anotaciones (AnnA) es una tupla:

$$\mathcal{N} = (Q, \Sigma, \Lambda, \Delta, I, F)$$

- $Q$  es un conjunto finito de estados.
- $\Sigma$  es el alfabeto de input.
- $\Lambda$  es un conjunto finito de etiquetas (Labels).
- $\Delta \subseteq Q \times (\Sigma \cup \Sigma \times \mathcal{N}) \times Q$  es la relación de transición.
- $I \subseteq Q$  es un conjunto de estados iniciales.
- $F \subseteq Q$  es el conjunto de estados finales (o aceptación).

Las transiciones  $(p, a, \ell, q)$  simbolizan que al leer la letra  $a$ , está letra **será anotada** con  $\ell$

# Ejecución de un AnnA (recordatorio)

Sea un AnnA  $\mathcal{N} = (Q, \Sigma, \Lambda, \Delta, I, F)$  y un documento  $d = a_0 \dots a_{n-1} \in \Sigma^*$ .

## Definiciones

Una **ejecución**  $\rho$  de  $\mathcal{A}$  sobre  $d$  es una secuencia  $\rho : p_0 \xrightarrow{t_0} p_1 \xrightarrow{t_1} \dots \xrightarrow{t_{n-1}} p_n$  tal que cumple todas las siguientes condiciones:

- $p_0 \in I$
- para todo  $i < n$ ,  $t_i$  es de la forma  $t_i = a_i$  o  $t_i = (a_i, \ell)$  para algún  $\ell \in \Lambda$
- para todo  $i < n$ ,  $(p_i, t_i, p_{i+1}) \in \Delta$ .

Se define la **anotación de  $\rho$**  como  $\text{ann}(\rho) = \text{ann}_0(t_0) \cdot \dots \cdot \text{ann}_{n-1}(t_{n-1})$ :

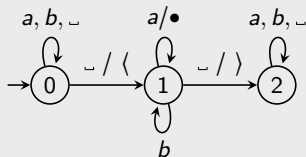
$$\text{ann}_i(t) = \begin{cases} (i, \ell) & t = (a, \ell) \\ \epsilon & t = a \end{cases}$$

Decimos que  $\rho$  es de **aceptación** ssi  $q_n \in F$ .

# Ejecución de un AnnA (recordatorio)

## Ejemplos de ejecuciones

$$d = \frac{a \quad a b b a \quad b a \quad b}{\begin{array}{cccccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \end{array}}$$



Algunas ejecuciones sobre  $d$ :

$$\rho_1 : p_0 \xrightarrow{a} p_0 \xrightarrow{\_ / <} p_1 \xrightarrow{a / \bullet} p_1 \xrightarrow{b} p_1 \xrightarrow{b} p_1 \xrightarrow{a / \bullet} p_1 \xrightarrow{\_ / >} p_2 \xrightarrow{b} p_2 \xrightarrow{a} p_2 \xrightarrow{\_} p_2 \xrightarrow{b} p_2$$

$$\rho_2 : p_0 \xrightarrow{a} p_0 \xrightarrow{\_} p_0 \xrightarrow{a} p_0 \xrightarrow{b} p_0 \xrightarrow{b} p_0 \xrightarrow{a} p_0 \xrightarrow{\_ / <} p_1 \xrightarrow{b} p_1 \xrightarrow{a / \bullet} p_1 \xrightarrow{\_ / >} p_2 \xrightarrow{b} p_2$$

$$\text{ann}(\rho_1) = (1, \langle) (2, \bullet) (5, \bullet) (6, \rangle) \qquad \text{ann}(\rho_2) = (6, \langle) (8, \bullet) (9, \rangle)$$

$$a \_ \dot{a} b b \dot{a} \_ b a \_ b$$

$$a \_ a b b a \_ b \dot{a} \_ b$$

# Output de un AnnA

Sea  $\mathcal{N} = (Q, \Sigma, \Lambda, \Delta, I, F)$  un autómata con anotaciones (AnnA).

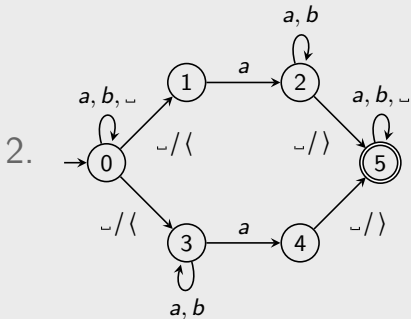
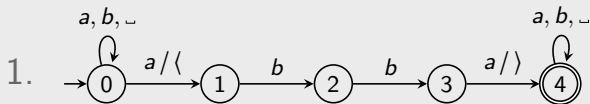
## Definición

Se define la función  $\llbracket \mathcal{N} \rrbracket$  tal que para todo documento  $d \in \Sigma^*$ :

$$\llbracket \mathcal{N} \rrbracket(d) = \{ \text{ann}(\rho) \mid \rho \text{ es una ejecución } \mathbf{aceptación} \text{ de } \mathcal{N} \text{ sobre } d \}$$

# Determinismo para autómatas con anotaciones

¿cuál de los dos AnnA diría que son deterministas?



# Determinismo para autómatas con anotaciones

Sea  $\mathcal{N} = (Q, \Sigma, \Lambda, \Delta, I, F)$  un autómata con anotaciones.

## Definición

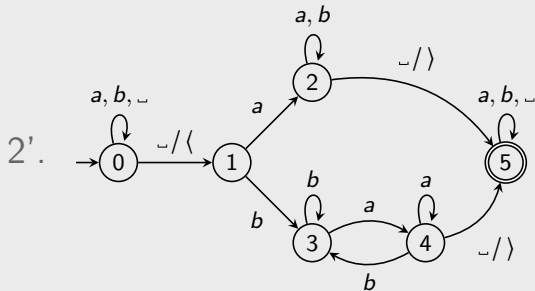
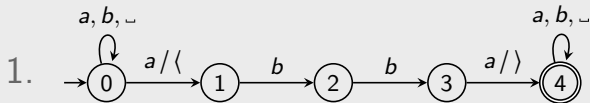
Decimos que  $\mathcal{N}$  es **Input-Output determinista** (I/O-determinista) ssi  $|I| = 1$  y para todo  $(p, t_1, q_1), (p, t_2, q_2) \in \Delta$ , si  $t_1 = t_2$ , entonces  $q_1 = q_2$ .

$\mathcal{N}$  es determinista al recibir  
el documento y una anotación **simultáneamente**

¿qué ventaja tiene un autómata I/O-determinista?

# Determinismo para autómatata con anotaciones

## Ejemplos





# Todo AnnA se puede I/O-determinizar

## Teorema

Para todo AnnA  $\mathcal{N} = (Q, \Sigma, \Lambda, \Delta, I, F)$ ,  
existe un AnnA I/O-determinista  $\mathcal{N}^{\text{det}}$  tal que  $\llbracket \mathcal{N} \rrbracket = \llbracket \mathcal{N}^{\text{det}} \rrbracket$ .

## Demostración

Considere la determinización de  $\mathcal{N}$  como:

$$\mathcal{N}^{\text{det}} = (2^Q, \Sigma, \Lambda, \Delta^{\text{det}}, q_0^{\text{det}}, F^{\text{det}})$$

- $2^Q = \{S \mid S \subseteq Q\}$  es el conjunto potencia de  $Q$ .
- $q_0^{\text{det}} = I$ .
- $\Delta^{\text{det}} : 2^Q \times (\Sigma \cup \Sigma \times \Lambda) \rightarrow 2^Q$  tal que para todo  $t \in \Sigma \cup (\Sigma \times \Lambda)$ :

$$\Delta^{\text{det}}(S, t) = \{q \in Q \mid \exists p \in S. (p, t, q) \in \Delta\}$$

- $F^{\text{det}} = \{S \subseteq Q \mid S \cap F \neq \emptyset\}$ .



# ¿cómo evaluamos una regex?

PROBLEMA: Evaluación de regex

INPUT: una regex  $R$  y  
un documento  $d$

OUTPUT: Enumerar todos los mappings en  $\llbracket R \rrbracket(d)$

1. Transformamos  $R$  a un vset autómata funcional  $\mathcal{A}_R$ .
2. Transformamos  $\mathcal{A}_R$  a un AnnA  $\mathcal{N}_R$ .
3. Determinizamos  $\mathcal{N}_R$  a un AnnA I/O determinista  $\mathcal{N}_R^{\text{det}}$ .
4. Calculamos los resultados  $\llbracket \mathcal{N}_R^{\text{det}} \rrbracket(d \cdot \#)$ .

¿cómo computamos todas las anotaciones de  $\llbracket \mathcal{N}_R^{\text{det}} \rrbracket(d \cdot \#)$ ?

# Outline

Enumeración con delay constante

Enumerable compact sets

Algoritmo de evaluación para AnnA

Algunas conclusiones

# Outline

**Enumeración con delay constante**

Enumerable compact sets

Algoritmo de evaluación para AnnA

Algunas conclusiones

# Problemas de enumeración

## Definición

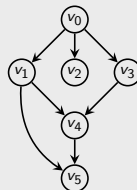
- Sea  $\mathcal{C}$  una clase de instancias (inputs) y  $\Omega$  un alfabeto (quizás infinito).
- Sea  $[\cdot] : \mathcal{C} \rightarrow 2^{\Omega^*}$  una función tal que para cada  $I$ ,  
 $[I] = \{out_1, out_2, \dots, out_k\}$  es un conjunto finito de outputs  $out_i \in \Omega^*$ .

## Ejemplo

$\mathcal{C}$ : la clase de todos los grafos  $G = (V, E)$  tal que  $G$  es acíclico (no tiene ciclos) y tiene un solo nodo raíz  $v_0$ .

$\Omega$ : conjunto de vértices.

$[G]$ : todos los caminos (secuencias) máximos desde  $v_0$ .



$$[G] = \{ v_0 v_1 v_5, v_0 v_1 v_4 v_5, \\ v_0 v_2, v_0 v_3 v_4 v_5 \}$$

# Problemas de enumeración

## Definición problema de enumeración

INPUT: Instancia  $I \in \mathcal{C}$ .

OUTPUT: Enumerar  $\llbracket I \rrbracket = \{out_1, out_2, \dots, out_k\}$ .

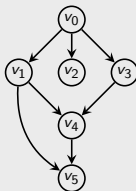
Con **enumerar** nos referimos a computar e imprimir en memoria:

$$\#out_1\#out_2\#\dots\#out_k\#$$

para algún orden de  $\llbracket I \rrbracket$  y donde  $\#$  es un símbolo separador especial.

## Ejemplo

INPUT:



OUTPUT:  $\# v_0 v_1 v_5 \# v_0 v_1 v_4 v_5 \# v_0 v_2 \# v_0 v_3 v_4 v_5 \#$

# Problemas de enumeración

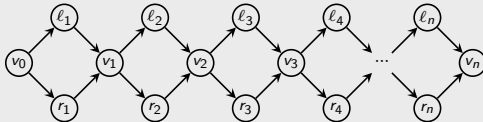
## Definición problema de enumeración

INPUT: Instancia  $I \in \mathcal{C}$ .

OUTPUT: Enumerar  $\llbracket I \rrbracket = \{out_1, out_2, \dots, out_k\}$ .

## Ejemplo problemático

INPUT:  $G$ :



■ ¿cuántos nodos tiene el grafo?

$n + 2n + 1$  nodos.

■ ¿cuántos caminos máximos tiene el grafo?

$2^n$  nodos.

¿cómo encontrar un algoritmo **eficiente** si el output puede ser exponencial?

# Algoritmos de enumeración con delay constante

## Definición problema de enumeración

INPUT: Instancia  $I \in \mathcal{C}$ .

OUTPUT: Enumerar  $\llbracket I \rrbracket = \{out_1, out_2, \dots, out_k\}$ .

Un algoritmo de enumeración para el problema anterior tiene **delay constante** si funciona en dos fases:

1. **Preprocesamiento**: el algoritmo recibe  $I$  y realiza algún cómputo.
2. **Enumeración**: terminada 1., el algoritmo imprime secuencialmente:

$$\begin{array}{ccccccccc} \# & out_1 & \# & out_2 & \# & \dots & \# & out_k & \# \\ \uparrow & & \uparrow & & \uparrow & & \uparrow & & \uparrow \\ 0 & & 1 & & 2 & & k-1 & & k \end{array}$$

donde existe  $C \in \mathbb{N}$  tal que para todo  $i \leq k$  el tiempo entre imprimir el  $(i-1)$ -ésimo y el  $i$ -ésimo símbolo  $\#$  es menor a  $C \cdot |out_i|$ .

El tiempo para entregar el siguiente output **solo depende de su tamaño**.



# Algoritmos de enumeración con delay constante

1. **Preprocesamiento:** el algoritmo recibe  $I$  y realiza algún cómputo.
2. **Enumeración:** terminada 1., el algoritmo imprime secuencialmente:

$\#out_1\#out_2\#\dots\#out_k\#$

donde existe  $C \in \mathbb{N}$  tal que para todo  $i \leq k$  el tiempo entre imprimir el  $(i-1)$ -ésimo y el  $i$ -ésimo símbolo  $\#$  es menor a  $C \cdot |out_i|$ .

**Function** enum-maxpaths ( $G = (V, E)$ )

Preprocesar  $G$  y encontrar  $v_0$

rec-maxpaths ( $G, v_0, v_0$ )

**print**  $\#$

**Function** rec-maxpaths ( $G, v, s$ )

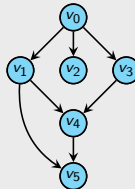
**if**  $\{u \mid (v, u) \in E\} = \emptyset$  **then**

**print**  $\# \cdot s$

**else**

**foreach**  $u \in \{u \mid (v, u) \in E\}$  **do**

rec-maxpaths ( $G, u, s \cdot u$ )



$\#v_0v_1v_5\#v_0v_1v_4v_5\#v_0v_2\#v_0v_3v_4v_5\#$

# Algoritmos de enumeración con delay constante

Un algoritmo de enumeración para el problema anterior tiene **delay constante** si funciona en dos fases:

1. **Preprocesamiento**: el algoritmo recibe  $I$  y realiza algún cómputo.
2. **Enumeración**: terminada 1., el algoritmo imprime secuencialmente:

$$\#out_1\#out_2\#\dots\#out_k\#$$

donde existe  $C \in \mathbb{N}$  tal que para todo  $i \leq k$  el tiempo entre imprimir el  $(i-1)$ -ésimo y el  $i$ -ésimo símbolo  $\#$  es menor a  $C \cdot |out_i|$ .

## Definición

Si el algoritmo de enumeración con delay constante toma tiempo  $\mathcal{O}(f(|I|))$  diremos que nuestro algoritmo **toma tiempo  $\mathcal{O}(f)$** .

## Ejemplo

Para el algoritmo de enumeración para caminos máximos, el preprocesamiento (ej. encontrar  $v_0$ ) tomará tiempo  $\mathcal{O}(|G|)$ .

Por lo tanto, el algoritmo de enumeración toma **tiempo lineal** en  $G$ .

# ¿cómo evaluamos una regex?

1. Transformamos  $R$  a un vset autómata funcional  $\mathcal{A}_R$ .
2. Transformamos  $\mathcal{A}_R$  a un AnnA  $\mathcal{N}_R$ .
3. Determinizamos  $\mathcal{N}_R$  a un AnnA I/O determinista  $\mathcal{N}_R^{\text{det}}$ .
4. Calculamos los resultados  $\llbracket \mathcal{N}_R^{\text{det}} \rrbracket(d \cdot \#)$ .

PROBLEMA: Evaluación de AnnA I/O-determinista

INPUT: un AnnA I/O-determinista  $\mathcal{N}$  y  
un documento  $d$

OUTPUT: Enumerar todas las anotaciones en  $\llbracket \mathcal{N} \rrbracket(d)$ .

Queremos una algoritmo de enumeración  
con **delay constante** y **tiempo**  $\mathcal{O}(|\mathcal{N}| \cdot |d|)$ .

# Outline

Enumeración con delay constante

**Enumerable compact sets**

Algoritmo de evaluación para AnnA

Algunas conclusiones

# Enumerable Compact Set

Sea  $\Omega$  un alfabeto (quizás infinito).

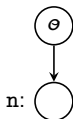
## Definición

Un **enumerable compact set** (ECS) es una estructura de datos compuesta por nodos y conexiones, donde los nodos pueden ser de 3 tipos:

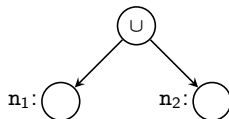
nodo  $\epsilon$



nodo  $\theta \in \Omega$



nodo  $\cup$



Cada nodo define un conjunto finito en  $\Omega^*$ , recursivamente:

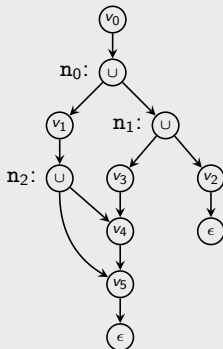
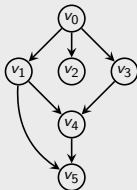
$$\llbracket \epsilon \rrbracket = \{\epsilon\}$$

$$\llbracket \theta \rrbracket = \llbracket n \rrbracket \cdot \{\theta\}$$

$$\llbracket \cup \rrbracket = \llbracket n_1 \rrbracket \cup \llbracket n_2 \rrbracket$$

# Enumerable Compact Set

## Ejemplo de una ECS y sus resultados



$$\llbracket v_4 \rrbracket = \{v_5 v_4\}$$

$$\llbracket n_2 \rrbracket = \{v_5, v_5 v_4\}$$

$$\llbracket v_1 \rrbracket = \{v_5 v_1, v_5 v_4 v_1\}$$

$$\llbracket n_1 \rrbracket = \{v_5 v_4 v_3, v_2\}$$

$$\llbracket n_0 \rrbracket = \{v_5 v_1, v_5 v_4 v_1, v_5 v_4 v_3, v_2\}$$

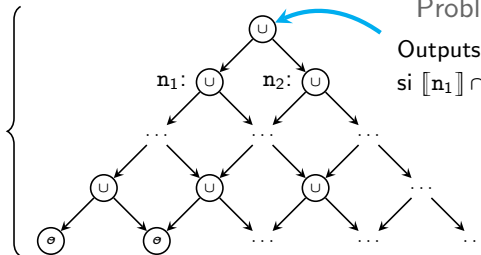
# Problema de enumeración para ECS

INPUT: un ECS  $\mathcal{E}$  y  
un nodo  $n$  de  $\mathcal{E}$

OUTPUT: Enumerar todos los resultados  $\llbracket n \rrbracket$ .

¿podemos **enumerar** todos los resultados con **delay constante**?

Problema 1.  
Delay dependerá  
de esta altura



Problema 2.

Outputs repetidos  
si  $\llbracket n_1 \rrbracket \cap \llbracket n_2 \rrbracket \neq \emptyset$

# Restricciones a Enumerable Compact Set

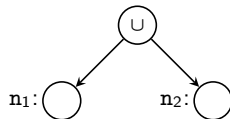
nodo  $\epsilon$



nodo  $\theta \in \Omega$



nodo  $\cup$



## Definiciones

0. Se define el odepth de un nodo recursivamente como:

$$\text{odepth}(\epsilon) = 0 \quad \text{odepth}(\theta) = 0 \quad \text{odepth}(\cup) = \text{odepth}(n_1) + 1$$

1. Un ECS  $\mathcal{E}$  es **k-acotado** si  $\text{odepth}(n) \leq k$  para todo nodo  $n$  en  $\mathcal{E}$ .
2. Un ECS  $\mathcal{E}$  es **sin repeticiones** si para todo nodo  $\cup$  (como arriba) en  $\mathcal{E}$  se cumple que  $\llbracket n_1 \rrbracket \cap \llbracket n_2 \rrbracket = \emptyset$ .



# Enumeración con delay constante para ECS

## Definiciones

0. Se define el odepth de un nodo recursivamente como:

$$\text{odepth}(\odot_{\epsilon}) = 0 \quad \text{odepth}(\odot_{\emptyset}) = 0 \quad \text{odepth}(\odot_{\cup}) = \text{odepth}(n_1) + 1$$

1. Un ECS  $\mathcal{E}$  es ***k*-acotado** si  $\text{odepth}(n) \leq k$  para todo nodo  $n$  en  $\mathcal{E}$ .
2. Un ECS  $\mathcal{E}$  es ***sin repeticiones*** si para todo nodo  $\odot_{\cup}$  (como arriba) en  $\mathcal{E}$  se cumple que  $\llbracket n_1 \rrbracket \cap \llbracket n_2 \rrbracket = \emptyset$ .

## Teorema

Para  $k \geq 1$  fijo, existe un algoritmo que, para todo ECS ***k*-acotado** y ***sin repeticiones***  $\mathcal{E}$  y nodo  $n$ , enumera  $\llbracket n \rrbracket$  con delay constante y tiempo  $\mathcal{O}(1)$ .

Demostración: ejercicio.

# Operaciones para un ECS

Deseamos soportar las siguientes **interfaz** sobre un ECS  $\mathcal{E}$ :

$$n_e := \mathcal{E}.\text{epsilon}$$

$$\llbracket n_e \rrbracket = \{\epsilon\}$$

$$n_x := \mathcal{E}.\text{extend}(n, \theta)$$

$$\llbracket n_x \rrbracket = \llbracket n \rrbracket \cdot \{\theta\}$$

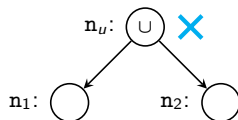
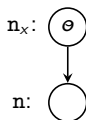
$$\text{t.q. } \theta \in \Omega$$

$$n_u := \mathcal{E}.\text{union}(n_1, n_2)$$

$$\llbracket n_u \rrbracket = \llbracket n_1 \rrbracket \cup \llbracket n_2 \rrbracket$$

$$\text{t.q. } \llbracket n_1 \rrbracket \cap \llbracket n_2 \rrbracket = \emptyset$$

¿cómo implementamos estas operaciones?



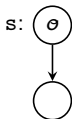
¿cómo soportamos estas operaciones manteniendo **k-acotado**?

# Nodos seguros

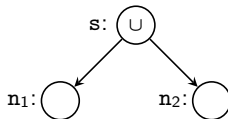
nodo  $\epsilon$



nodo  $\theta \in \Omega$



nodo  $\cup$

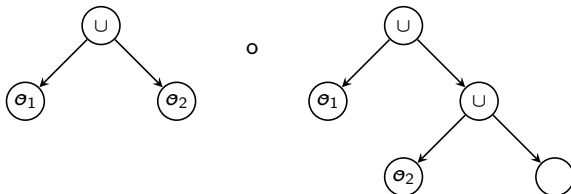


## Definición

Decimos que un nodo  $s$  es **seguro** (safe) si cumple que:

- $s$  es nodo  $\epsilon$  o nodo  $\theta$ , o
- $s$  es nodo  $\cup$ ,  $\text{odepth}(s) \leq 1$  y  $\text{odepth}(n_2) \leq 1$ .

En otras palabras, si  $s$  es nodo- $\cup$  y seguro, entonces  $s$  es de dos formas:



# Operaciones para un ECS sobre nodos seguros

Para nodos **seguros**  $s$ ,  $s_1$ ,  $s_2$  soportamos las operaciones sobre un ECS  $\mathcal{E}$ :

$$s_e := \mathcal{E}.\text{epsilon}$$

$$s_x := \mathcal{E}.\text{extend}(s, \theta)$$

$$s_u := \mathcal{E}.\text{union}(s_1, s_2)$$

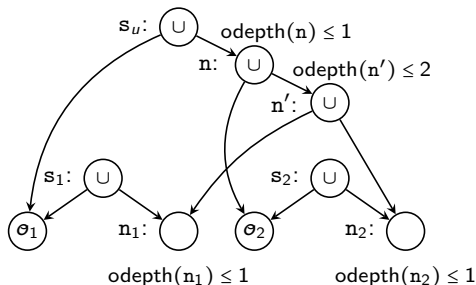
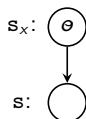
$$\llbracket s_e \rrbracket = \{\epsilon\}$$

$$\llbracket s_x \rrbracket = \llbracket s \rrbracket \cdot \{\theta\}$$

$$\llbracket s_u \rrbracket = \llbracket s_1 \rrbracket \cup \llbracket s_2 \rrbracket$$

$$\text{t.q. } \theta \in \Omega$$

$$\text{t.q. } \llbracket s_1 \rrbracket \cap \llbracket s_2 \rrbracket = \emptyset$$



Si  $\mathcal{E}$  es 2-**acotado**, entonces el resultado también es 2-**acotado**.

# Resumen sobre ECS

## Teorema

Para un ECS **2-acotado** y **sin repetidos**  $\mathcal{E}$ , y un nodo cualquiera  $n$  de  $\mathcal{E}$ :

1.  $\llbracket n \rrbracket$  se puede enumerar con **delay constante** y tiempo  $\mathcal{O}(1)$ .
2. Las operaciones **epsilon**, **extend** y **union** **sobre nodos seguros** mantienen la propiedad de **2-acotado** y **retorna nodos seguros**.
3. Los nodos input de las operaciones **mantienen sus resultados**.
4. Cada operación toma **tiempo constante**  $\mathcal{O}(1)$ .

Usaremos nuestra estructura ECS para la evaluación de AnnA

# Outline

Enumeración con delay constante

Enumerable compact sets

**Algoritmo de evaluación para AnnA**

Algunas conclusiones

# Algoritmo de evaluación con delay constante

Sea  $\mathcal{N} = (Q, \Sigma, \Lambda, \Delta, q_0, F)$  un **AnnA I/O-determinista**,  $d = a_0 \dots a_{n-1}$  un **documento**,  $\mathcal{E}$  un **ECS**, y  $S$  un **arreglo de nodos** indexados por  $Q$ .

**Function** eval-delayconstante( $\mathcal{A}, d$ )

$S \leftarrow [\text{null}]$

$S[q_0] \leftarrow \mathcal{E}.\text{epsilon}$

**for**  $i = 0$  **to**  $n - 1$  **do**

$S_{\text{old}} \leftarrow S$

$S \leftarrow [\text{null}]$

**foreach**  $(p, a_i, q) \in \Delta$  **such that**  $S_{\text{old}}[p] \neq \text{null}$  **do**

**if**  $S[q] = \text{null}$  **then**  $S[q] \leftarrow S_{\text{old}}[p]$

**else**  $S[q] \leftarrow \mathcal{E}.\text{union}(S[q], S_{\text{old}}[p])$

**foreach**  $(p, a_i, \ell, q) \in \Delta$  **such that**  $S_{\text{old}}[p] \neq \text{null}$  **do**

$n \leftarrow \mathcal{E}.\text{extend}(S_{\text{old}}[p], (i, \ell))$

**if**  $S[q] = \text{null}$  **then**  $S[q] \leftarrow n$

**else**  $S[q] \leftarrow \mathcal{E}.\text{union}(S[q], n)$

**foreach**  $q \in \{p \in Q \mid S[p] \neq \text{null}\} \cap F$  **do**

**enumerate**( $\mathcal{E}, S[q]$ )

# Análisis del algoritmo

- $\mathcal{N}$  es I/O-determinista y, por lo tanto,  $\mathcal{E}$  es **sin duplicados**.
- Cada operación en  $\mathcal{E}$  toma tiempo constante.
- Los resultados en  $\mathcal{E}$  se pueden enumerar con delay constante.

...entonces, el algoritmo enumera  $[[\mathcal{N}]](d)$   
con delay constante y en tiempo  $\mathcal{O}(|\mathcal{N}| \cdot |d|)$ .



# Outline

Enumeración con delay constante

Enumerable compact sets

Algoritmo de evaluación para AnnA

**Algunas conclusiones**

# Sobre los resultados de extracción de información

Son un **resumen/ideas** de los siguientes artículos:

Maturana, Riveros, Vrgoc

*"Document Spanners for Extracting Incomplete Information."* PODS 2018.

Florenzano, Riveros, Ugarte, Vansummeren, Vrgoc

*"Constant Delay Algorithms for Regular Document Spanners."* PODS 2018.

Grez, Riveros, Ugarte

*"A Formal Framework for Complex Event Processing."* ICDT 2019.

Muñoz, Riveros

*"Streaming Enumeration on Nested Documents."* ICDT 2022.

Amarilli, Jachiet, Muñoz, Riveros

*"Efficient Enumeration for Annotated Grammars."* PODS 2022.

Bucchi, Grez, Quintana, Riveros, Vansummeren

*"CORE: a COMplex event Recognition Engine."* VLDB 2022.

+ artículos de otros grupos + investigación en curso.

Query

```
(^|\n)!x{[A-Z][a-z]{4,}} !y{([A-Z][a-z ]+)}($|\n)
```

 RUN

Text

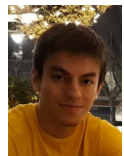
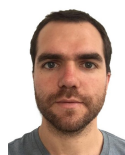
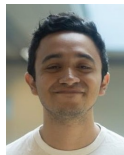
```
1 Nicolas Van Sint Jan
2 Vicente Calisto
3 Marjorie Bascunan
4 Oscar Carcamo
5 Cristian Riveros
6 Domagoj Vrgoc
```

 IMPORT FILE

Matches

id	x	y
0	Domagoj	Vrgoc
1	Cristian	Riveros
2	Oscar	Carcamo
3	Marjorie	Bascunan
4	Vicente	Calisto
5	Nicolas	Van_Sint_Jan

# Desarrollo de investigación / REmatch



(Faltan varios)

# Sobre mis temas de investigación

## Manejo de datos:

- Big data.
- Datos streaming.
- Extracción de información.

## Grafos de datos:

- Web semántica.
- Base de datos de grafos.
- Centralidad de datos.

## Lógica / Lenguajes formales:

- Teoría de modelos finitos.
- Teoría de automatas.

## Teoría de la Computación:

- Complejidad computacional.
- Algoritmos aleatorios.

# Proyectos de implementación

1. REmatch

**Motor de extracción de información**

2. CORE

**Base de datos streaming**

3. MilleniumDB

**Base de datos de grafos**

Están **invitados a colaborar** en cualquiera de estos proyectos

(solo escribanme y pregunten)

FIN :(

Gracias!