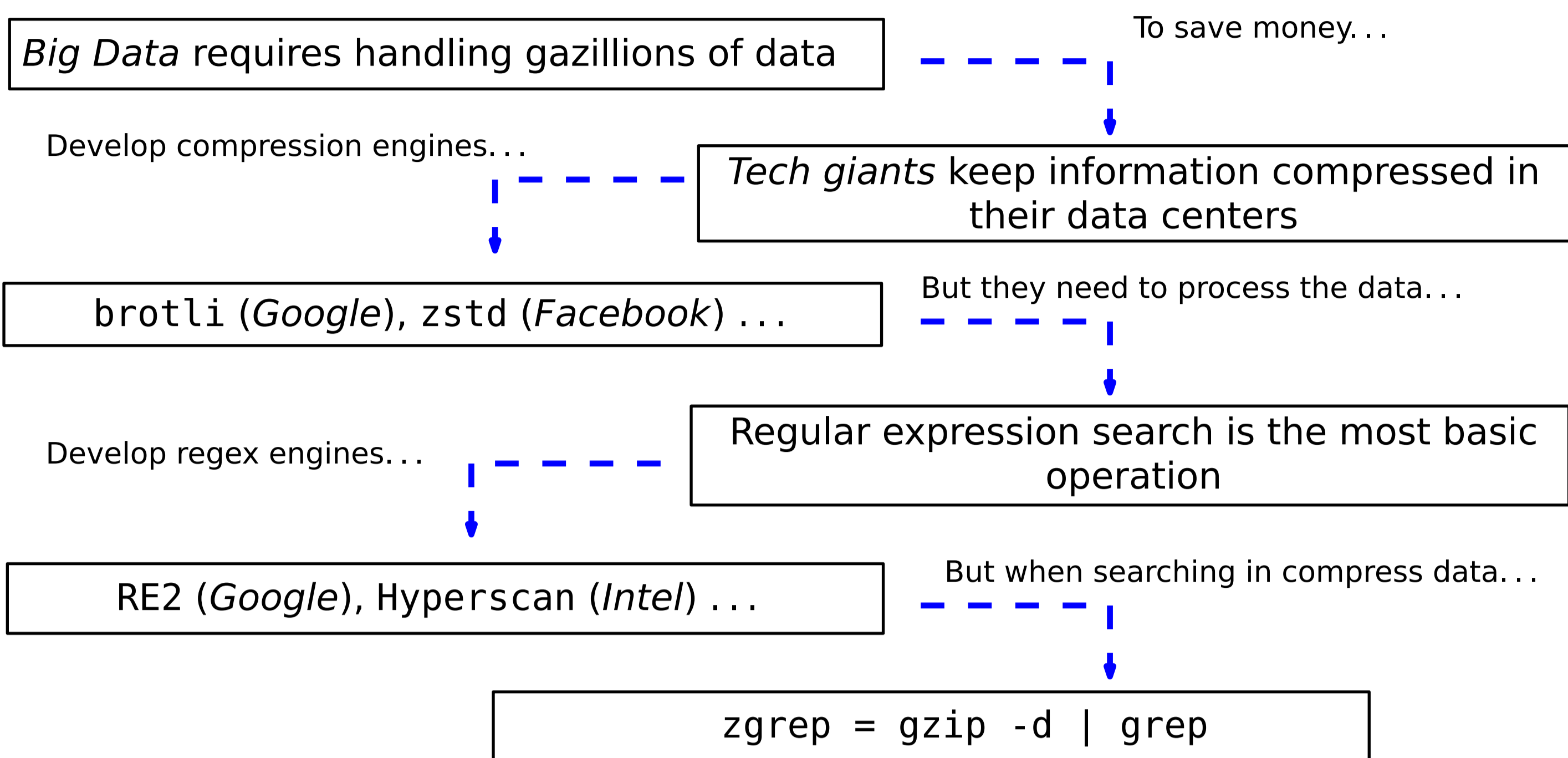


# Regular Expression Search on Compressed Text without the need for decompression

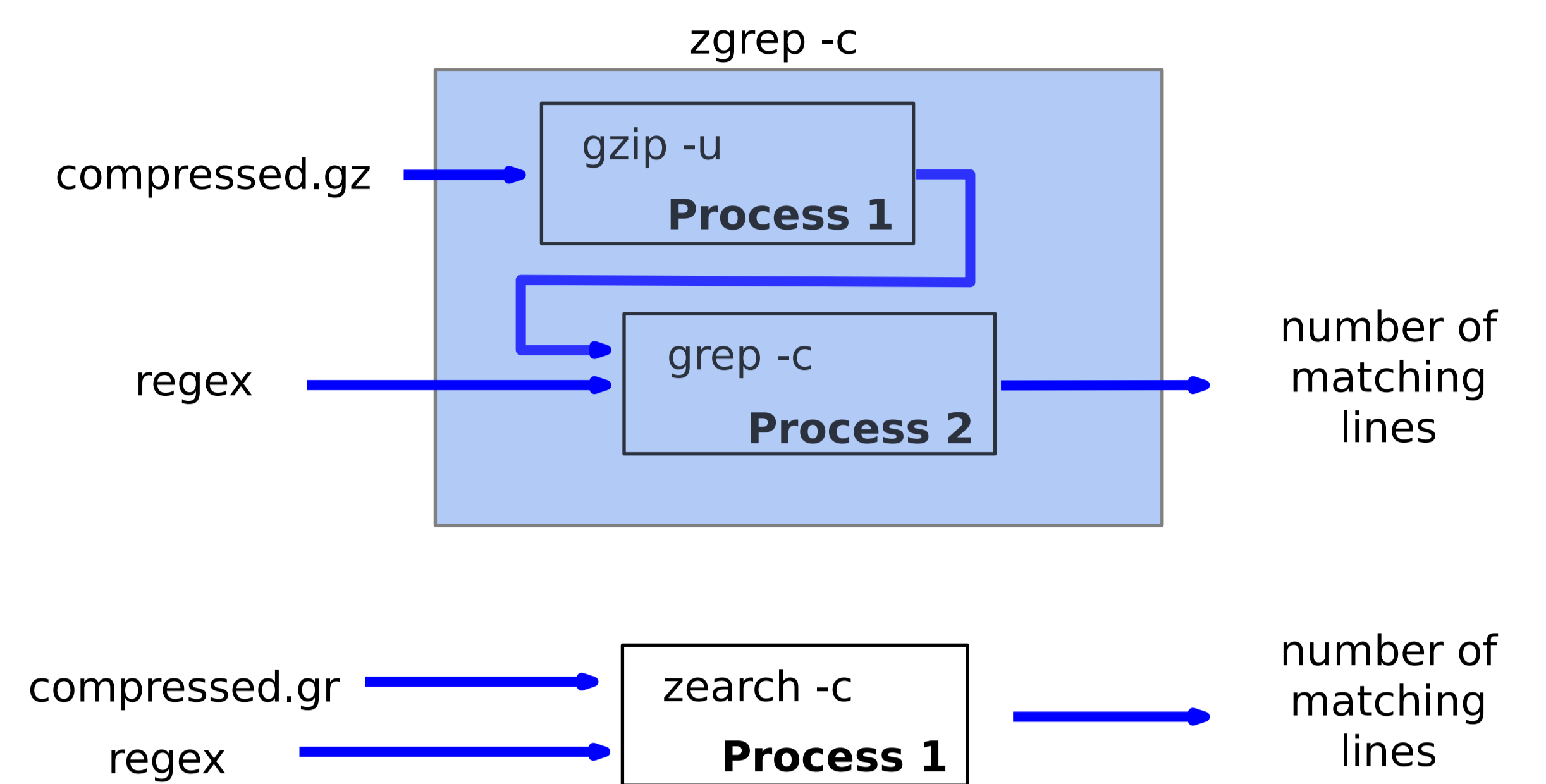
Javier Esparza, Technical University of Munich  
Pedro Valero & Pierre Ganty, IMDEA Software Institute

<https://pevalme.github.io>

## Problem



## State of the art vs zearch



## Overview of zearch

### Setup:

Data is compressed with a **grammar-based** algorithm (LZ78, LZW, Sequitur, Repair...)  
Regular Expression is used to build a **non-deterministic** automata without  $\epsilon$ -transitions

### Idea:

- Use redundancies to speed up the search.

Each rule is processed *only* once  
Processing  $n$  rules may cover  $2^n$  characters:  
 $\{X_0 \rightarrow X_1X_1\} \cup \{X_i \rightarrow X_{i+1}X_{i+1} \mid i = 1 \dots n\}$

- Saturation construction

- $(p) \xrightarrow{a} (q)$  NFA moves from  $p$  to  $q$  when reading  $a$
- $(p) \xrightarrow{X} (q)$  NFA moves from  $p$  to  $q$  when reading  $\mathcal{L}(X)$

### Algorithm:

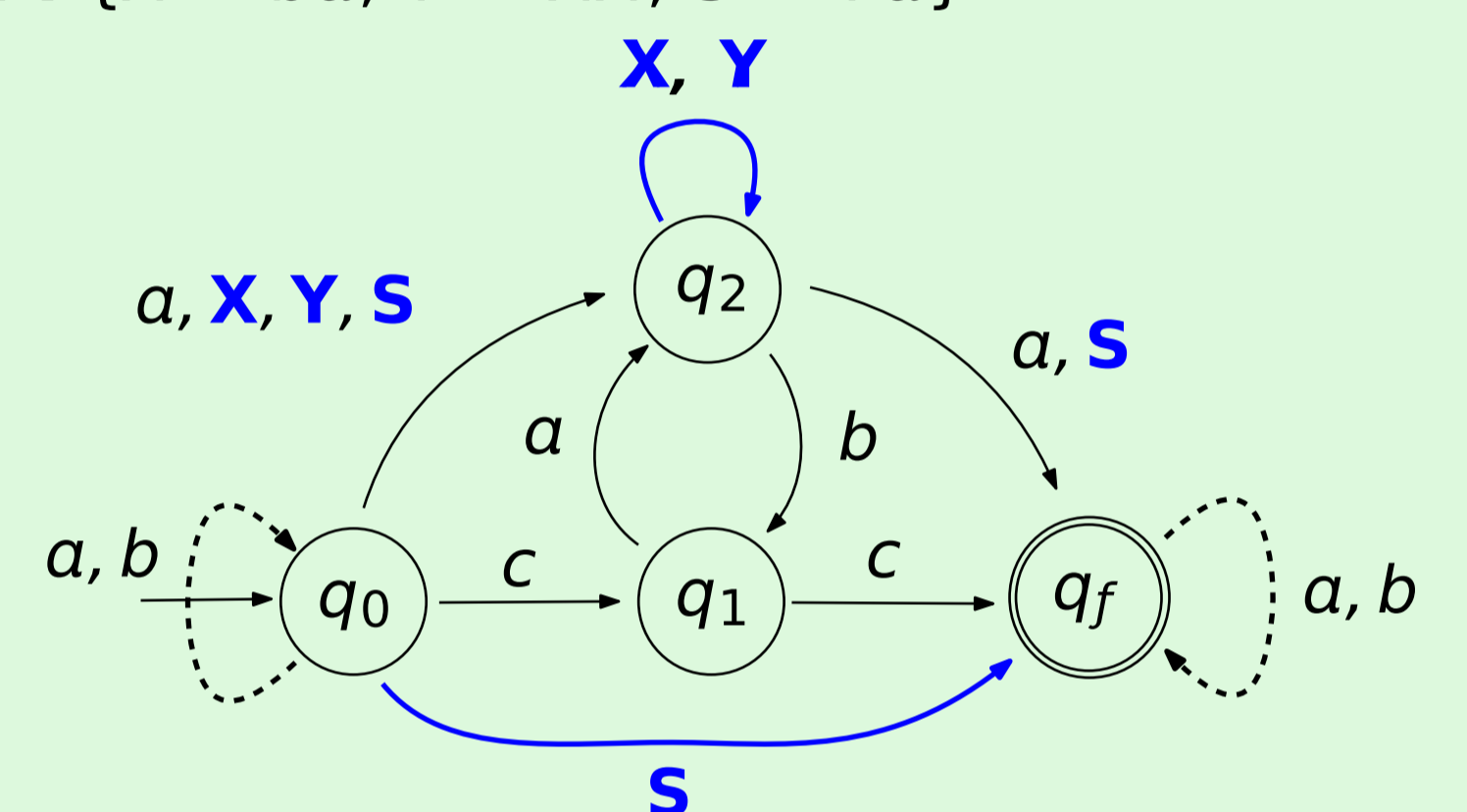
**Input:** SLP  $G = (\Sigma, V, \{X_1 \rightarrow \alpha_1\beta_1, \dots, X_n \rightarrow \alpha_n\beta_n\})$   
Automata  $A = (\Sigma, Q, q_0, F, \delta)$

**Output:** Saturated automata  $A^s = (\Sigma, Q, q_0, F, \delta^s)$

**For each**  $(X_i \rightarrow \alpha_i\beta_i) \in \{X_1 \rightarrow \alpha_1\beta_1, \dots, X_n \rightarrow \alpha_n\beta_n\}$   
**For each**  $(q_1, \alpha, q') \in \delta$   
**For each**  $(q', \beta, q_2) \in \delta$   
 $\delta := \delta \cup \{(q_1, X, q_2)\}$   
 $INCR\_COUNT(X, q')$

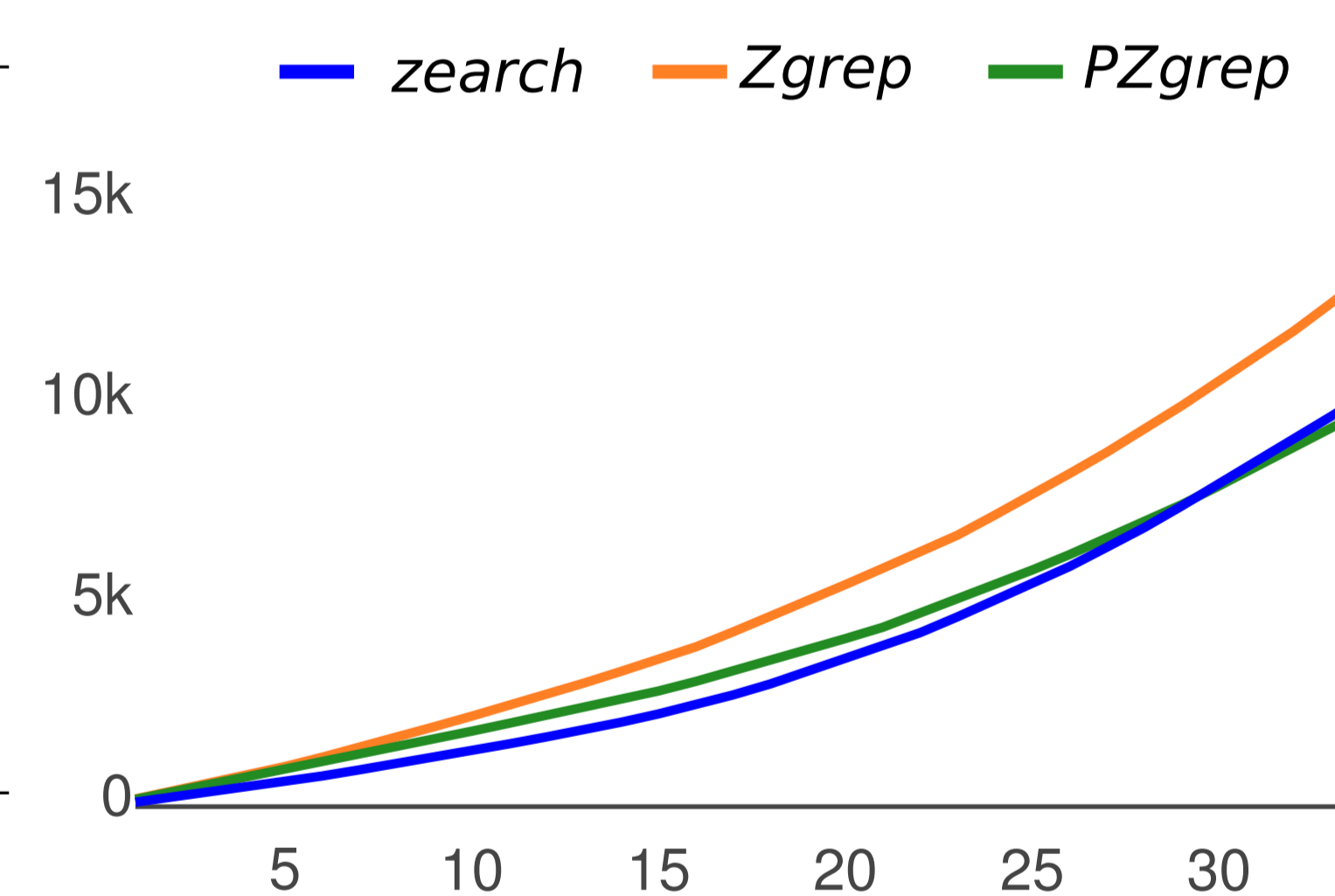
### Example:

SLP:  $\{X \rightarrow ba; Y \rightarrow XX; S \rightarrow Ya\}$



## Experimental Results

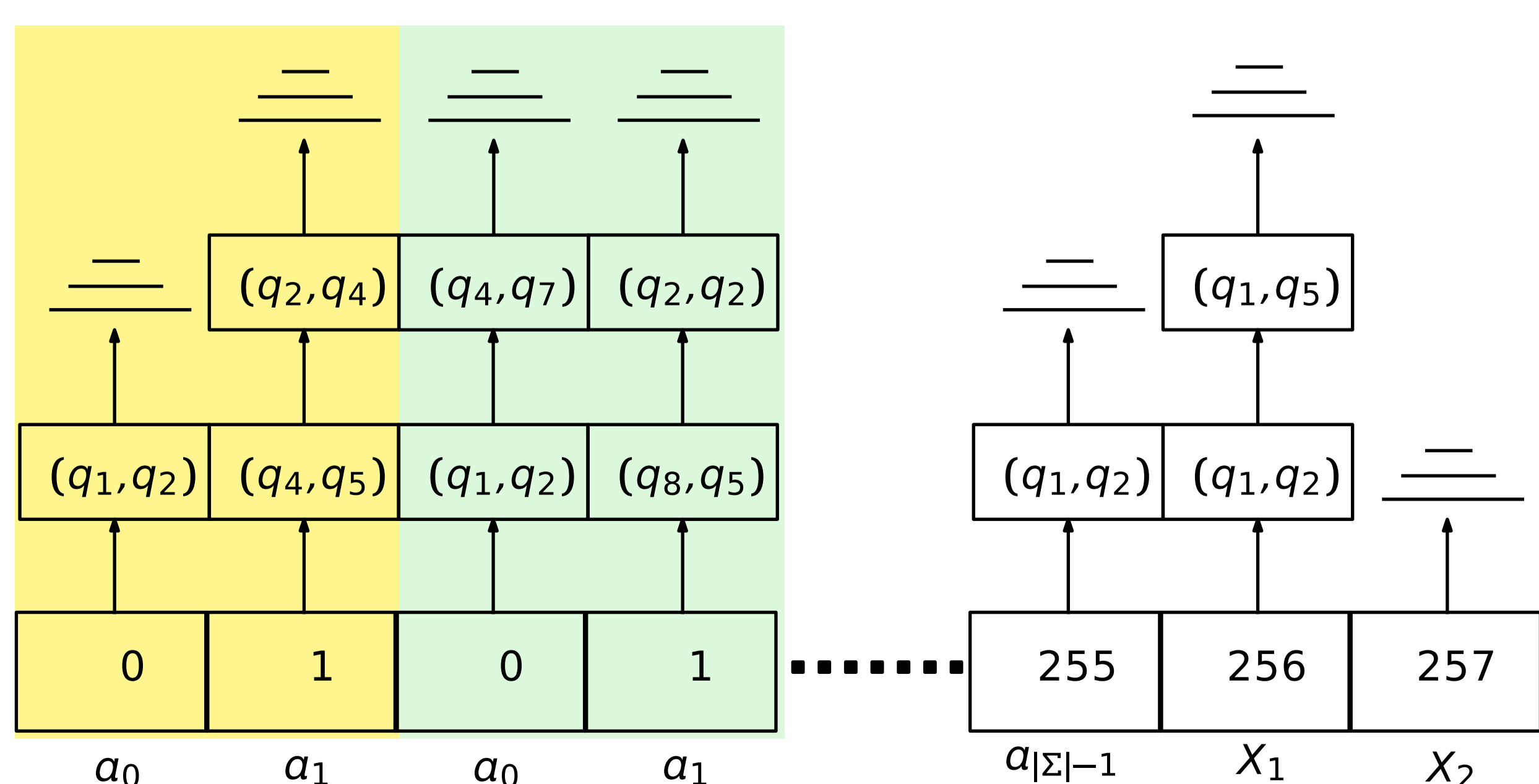
|                   | Log<br>100MB → 7 MB |       |        | JSON<br>100MB → 8.8 MB |       |            | Subtítulos<br>100MB → 14 MB |       |            | CSV<br>100MB → 27 MB |            |            |
|-------------------|---------------------|-------|--------|------------------------|-------|------------|-----------------------------|-------|------------|----------------------|------------|------------|
|                   | zearch              | Zgrep | PZgrep | zearch                 | Zgrep | PZgrep     | zearch                      | Zgrep | PZgrep     | zearch               | Zgrep      | PZgrep     |
| "pedro"           | <b>110</b>          | 261   | 212    | <b>152</b>             | 280   | 211        | <b>242</b>                  | 305   | <b>236</b> | 505                  | 417        | <b>377</b> |
| "."               | <b>130</b>          | 227   | 197    | <b>158</b>             | 185   | 194        | <b>283</b>                  | 367   | <b>273</b> | 595                  | 529        | <b>403</b> |
| "I .* you"        | <b>130</b>          | 272   | 199    | <b>210</b>             | 290   | 213        | <b>391</b>                  | 384   | <b>277</b> | 541                  | 373        | <b>351</b> |
| "[a-z]{4}"        | <b>161</b>          | 270   | 206    | <b>278</b>             | 188   | <b>186</b> | <b>407</b>                  | 494   | <b>354</b> | 694                  | <b>621</b> | <b>479</b> |
| " [a-z]{4} "      | <b>136</b>          | 409   | 262    | <b>177</b>             | 229   | 212        | <b>398</b>                  | 695   | 519        | <b>568</b>           | 712        | <b>473</b> |
| "[a-z]*[a-z]{4}"  | <b>172</b>          | 266   | 205    | 324                    | 191   | <b>188</b> | 481                         | 521   | <b>399</b> | 755                  | 644        | <b>462</b> |
| "[a-z]*[a-z]{6}"  | <b>192</b>          | 302   | 222    | 396                    | 194   | <b>195</b> | <b>539</b>                  | 596   | <b>436</b> | 819                  | <b>650</b> | <b>480</b> |
| "\d5\d0\d4\d5\d " | <b>123</b>          | 416   | 269    | <b>153</b>             | 272   | 205        | <b>309</b>                  | 705   | 527        | 608                  | 724        | <b>510</b> |
| "([a-z]{5} )+"    | <b>134</b>          | 410   | 262    | <b>206</b>             | 235   | 210        | <b>614</b>                  | 706   | <b>531</b> | 575                  | 727        | <b>512</b> |
| "([a-z]{5} +){5}" | <b>174</b>          | 528   | 385    | <b>323</b>             | 545   | 404        | <b>652</b>                  | 751   | <b>583</b> | 744                  | 748        | <b>468</b> |
| <b>Average</b>    | <b>146</b>          | 336   | 241    | <b>237</b>             | 260   | <b>221</b> | <b>431</b>                  | 552   | <b>413</b> | 640                  | <b>614</b> | <b>451</b> |



Better compression  $\Rightarrow$  More repetitions in the data  $\Rightarrow$  Less grammar rules to be processed  $\Rightarrow$  zearch is faster

## Data Structure

$X_1 \rightarrow a_0a_1$     $X_2 \rightarrow a_3a_4$     $X_3 \rightarrow X_2X_1$



$X \rightarrow \sigma, Y \rightarrow \rho$  are independent iff  $X \notin \rho \wedge Y \notin \sigma$   
Independent rules can be processed simultaneously  
Conceptually easy parallelization

List of edges labeled by the variable

Array with an element per variable of the grammar