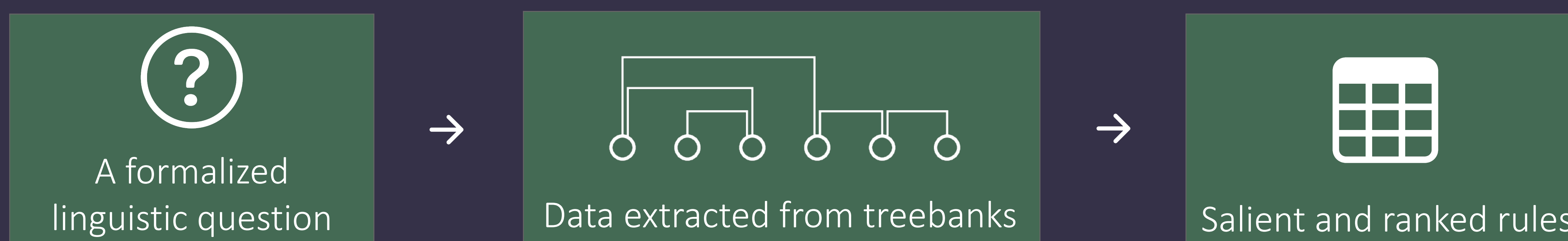


# Sparse Logistic Regression with High-order Features for Automatic Grammar Rule Extraction from Treebanks



Santiago Herrera, Caio Corro and Sylvain Kahane

Graphical abstract

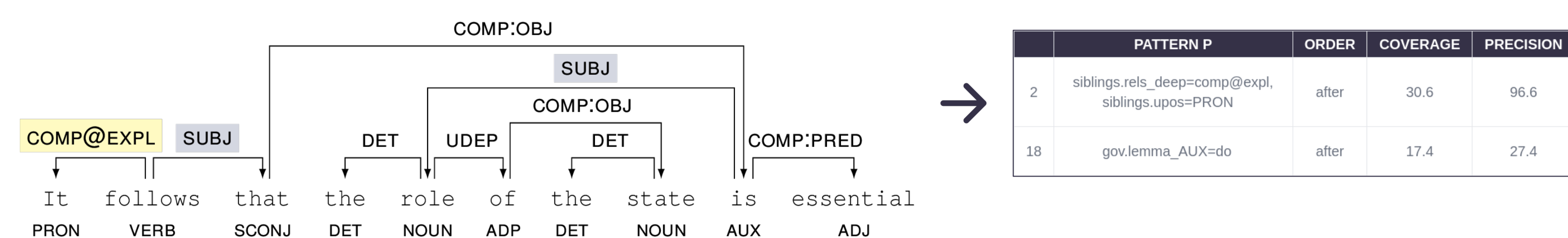


Expressive and quantitative fine-grained grammar rules for word-order and agreement

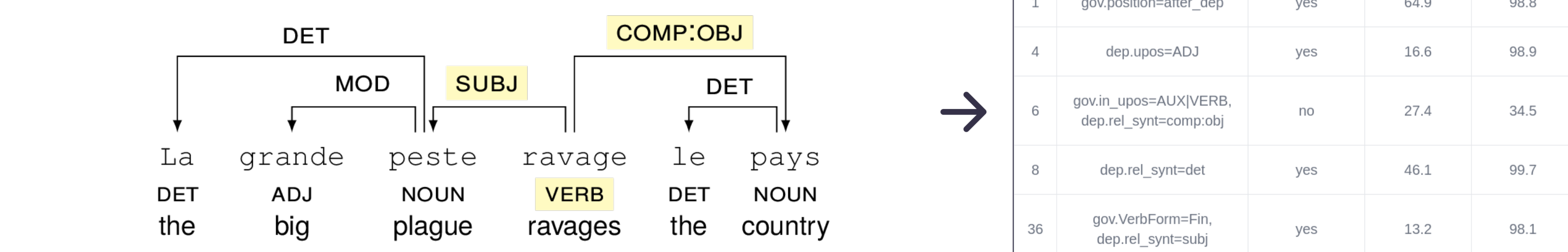
- Formalization of a syntactic rule
- ML method to extract and rank the rules
- An easy-to-interpret hierarchy of results

## 1 Grammatical descriptions

When the **subject follows the verb** in English?



When there is **number agreement** in French?



## 2 Definition of a syntactic rule

$$S \implies (P \xrightarrow{\alpha\%} Q).$$

e.g. Given all subjects (S), these are inverted (Q) when there is an expletive complement (P) in 96% of cases ( $\alpha$ ).

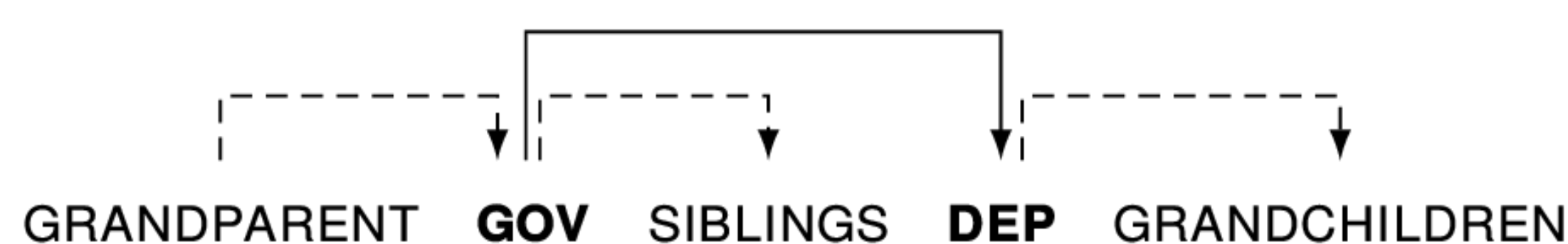
- Probabilistic
- Potentially overlapping
- +/- Fine-grained

## 3 Data and features

Work in 3 languages: **French, Spanish and Wolof**

Dependency treebanks : Universal Dependencies (UD) and Surface-Syntactic UD

The **P patterns** are filled in by the model using linguistic information within the following search space:



## 4 Rule extraction

**Sparse Logistic Regression**

$$P(\text{"number agreement"}|x) = \sigma(a^\top x + b) \text{ where } \sigma(w) = \frac{\exp w}{1 + \exp w}$$

**Training Problem**

$$\min_{a \in \mathbb{R}^F, b \in \mathbb{R}} \frac{1}{|D|} \sum_{(x,y) \in D} \ell(a^\top x + b; y) + \lambda r(a)$$



- Negative likelihood loss
- L1 Norm
- Ranks through the regularization parameter
- Sensible to less pronounced shift distributions

+ precision  
+ coverage/recall  
+ statistical tests

## R Results

More grammar rules with more expressiveness

An (almost) hyper-parameter free method to extract and rank rules

A better adapted tool for rule mining

