

Learning from Noisy Data

Learners use developing grammatical knowledge to parse & learn from their data
 ► How do they generalize accurately from **immature representations** of input?

Case Study: Word Order

Infants acquire their language's basic word order from input containing a mixture of canonical and non-canonical sentence types [1-4]

- (1) You're holding a toy.
- (2) What are you holding?
- (3) That's the dog we like.
- (4) You're being hugged.

		SVO?	
	English		French
	0.36 NP V	0.48 NP V	
	0.20 V	0.21 NP V NP	
	0.20 NP V NP	0.13 V	
	0.17 V NP	0.05 NP NP V	
	0.04 NP V NP NP	0.03 NP V NP NP	
	0.03 V NP NP	0.03 V NP	

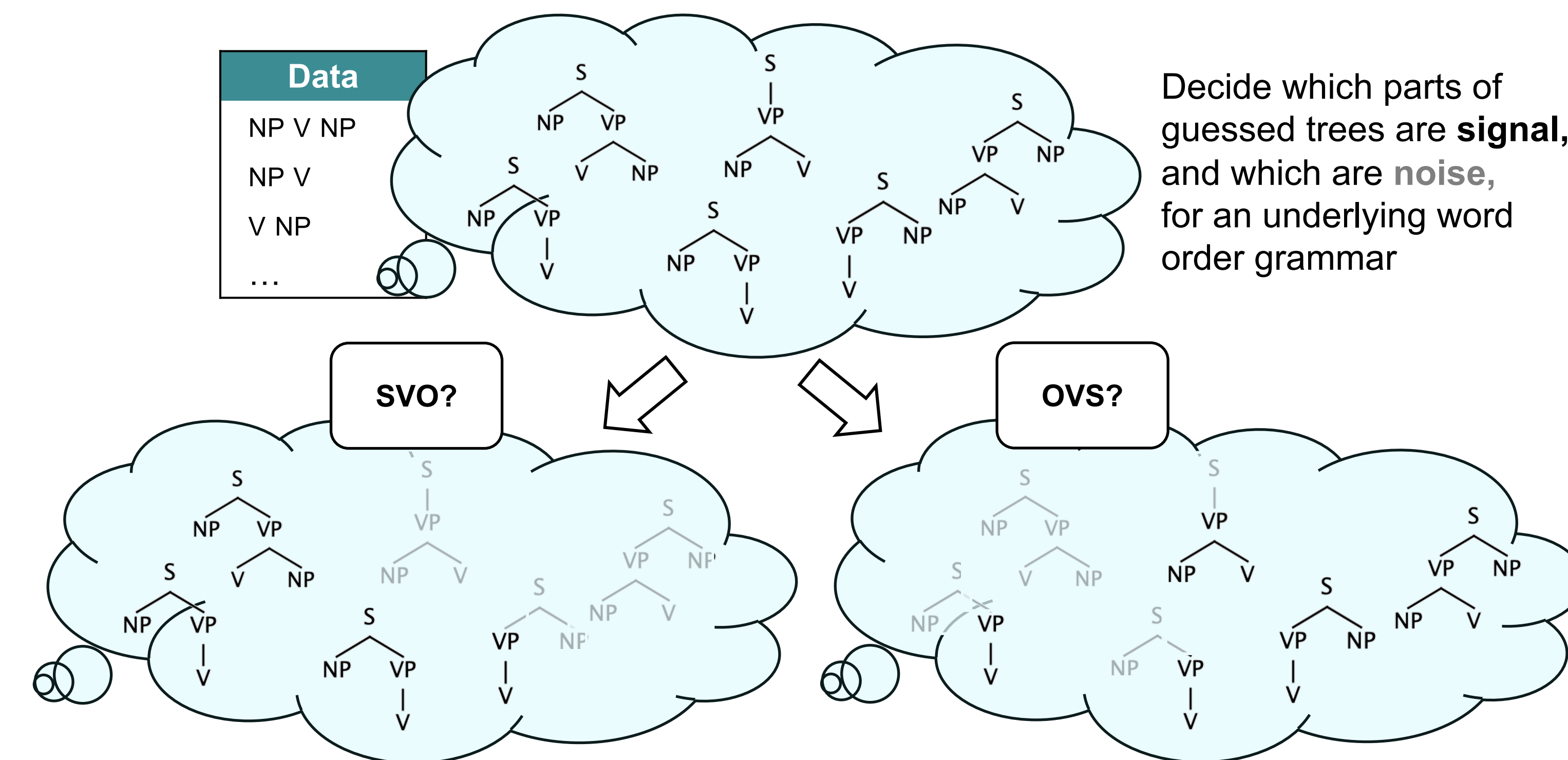
Fig. 1 Most frequent string types in Eve and Lyon CHILDES corpora [9-10]; NPs and Vs imperfectly identified from functional cues [11-13]

Distributions in child-directed speech are potentially misleading

- How do children avoid being misled by "noise" from non-canonical clause types? [4-8]

What does Filtering Look Like?

From strings of NPs and Vs, make a noisy guess about underlying tree structure



- It is possible to **learn** how to divide up the data into signal vs. noise, without knowing ahead of time how much noise there is, or what its properties are

Model Comparisons

"Fully-Flexible" Learner

No 4-way choice of a canonical word order grammar: all rules possible with some probability [18]

- Collapses distinction in our model between rules for canonical and non-canonical structures
- Learning canonical word order means identifying that some rules have probabilities near zero

Two variants: with and without an explicit bias to **regularize** (push probabilities towards zero/one) [15-17]

- Learner without bias to regularize infers distributions that mirror its noisy data
- Learner with bias to regularize gives high probability to non-target word orders

- Useful to have a hypothesis space with restrictive grammatical options

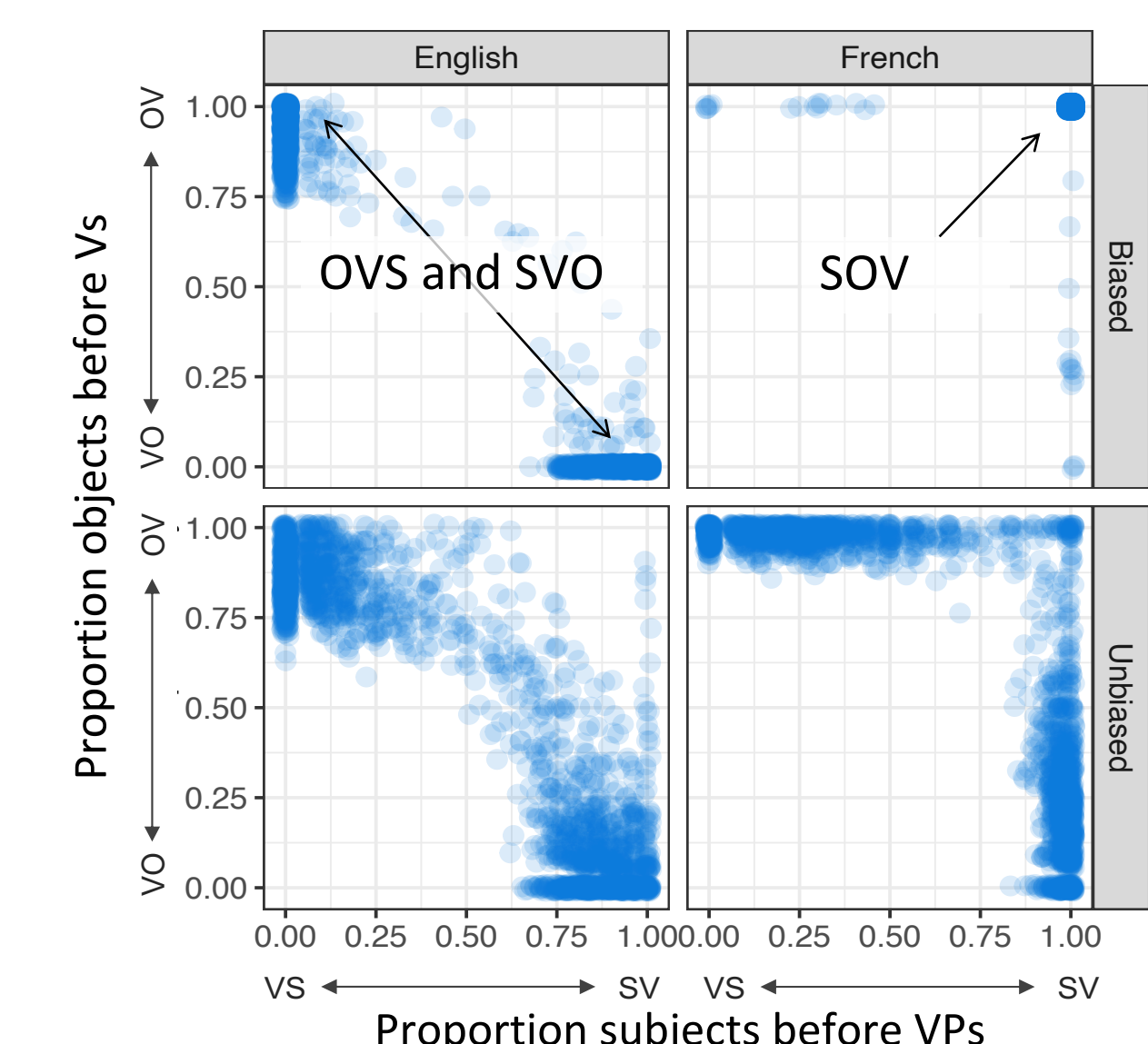


Fig. 3 Posterior distribution over subject and object positions in trees, fully-flexible learner

A Data-Coverage Heuristic

Simpler alternative: select the grammar that covers the most data

- E.g., core rules of SVO grammar generate 56% of English data, more than any other grammar

Comparison: version of our learner with an 8-way choice among grammars

- 4 options from original model that fix both subject and object position
- 4 less restrictive options that only fix one argument position, and allow the other to vary

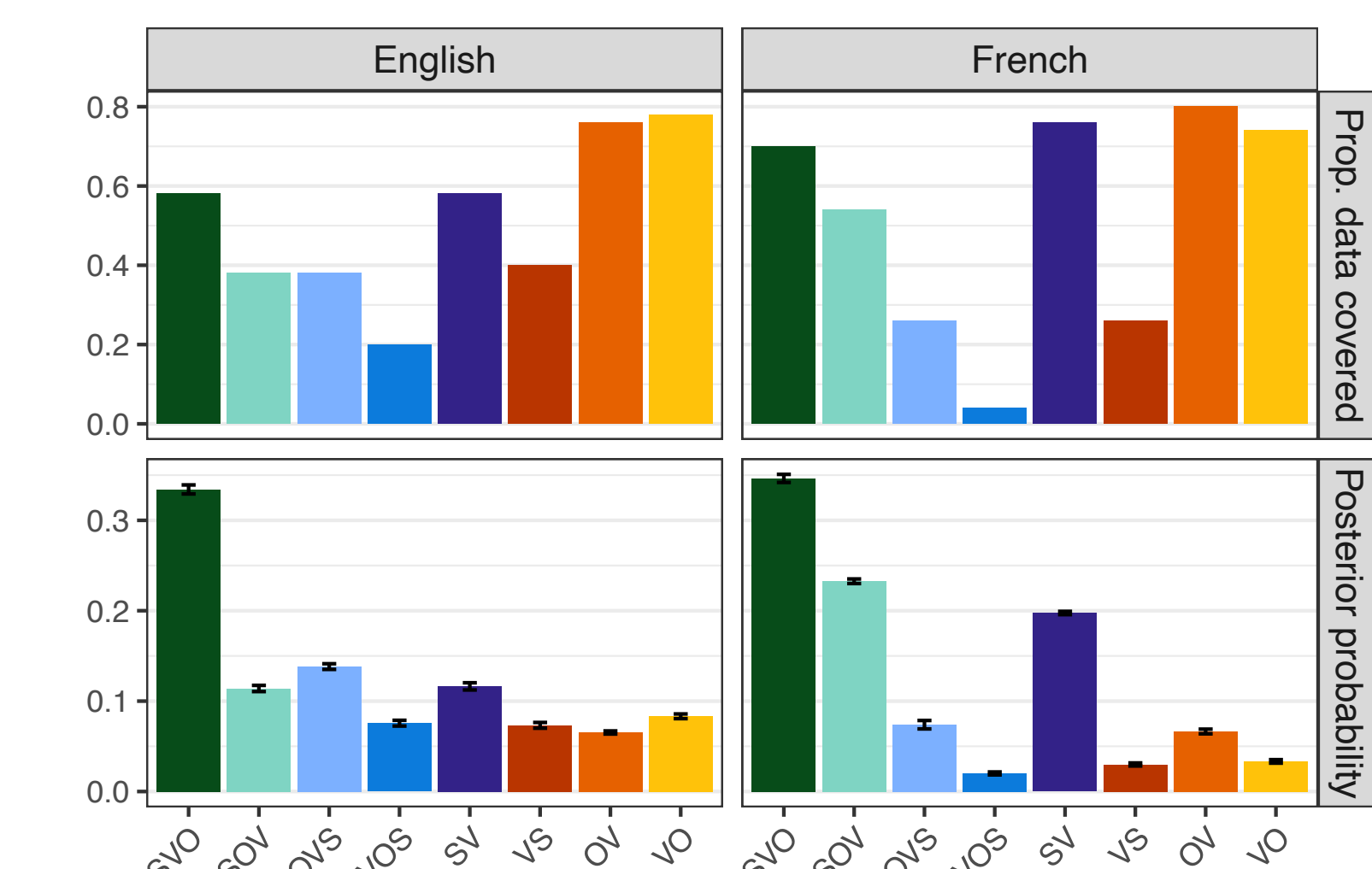


Fig. 4 Eight-way hypothesis space: proportion data coverage vs. model's posterior distribution

Our learner still successfully assigns SVO highest posterior probability in both languages, even though more flexible grammars cover more of the data

- Emergent preference for most restrictive hypothesis that fits the data

Proposal: Input Filtering

Two Possible Solutions

Regularization: explicit numerical bias against encoding full variability of data: prefer hypotheses that are heavily skewed [14-16]

Filtering: expect that data are a noisy realization of a deterministic underlying system, and learn to separate signal from noise [8]

- **Finding: filtering works better** in this learning domain

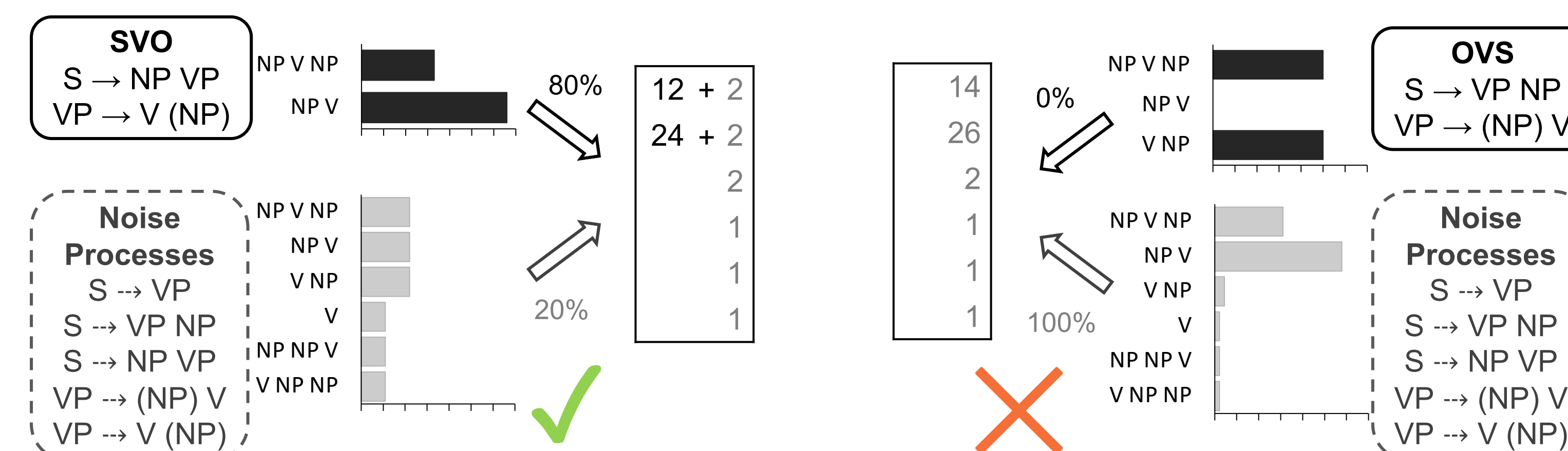
Toy Example

How might these data have arisen partially from a word order grammar distribution, and partially from the noise distribution?

NP V NP	14
NP V	26
V NP	2
V	1
NP NP V	1
V NP NP	1

SVO? OVS?
SOV? VOS?

Two solutions (of many):



- Costly to analyze too much of the data as noise: too many degrees of freedom
- Simpler solution: attribute skewed data to restrictive word order grammar whenever possible

Results: Child-Directed Speech

Simulations on 50-sentence datasets of NP-V strings, sampled from corpora of child-directed English and French (Fig. 1)

- Learner successfully assigns SVO highest posterior probability in both languages
- Even though data cannot be produced by any single word order grammar, without noise

- Filter works, and filter can be learned from distributions in the data

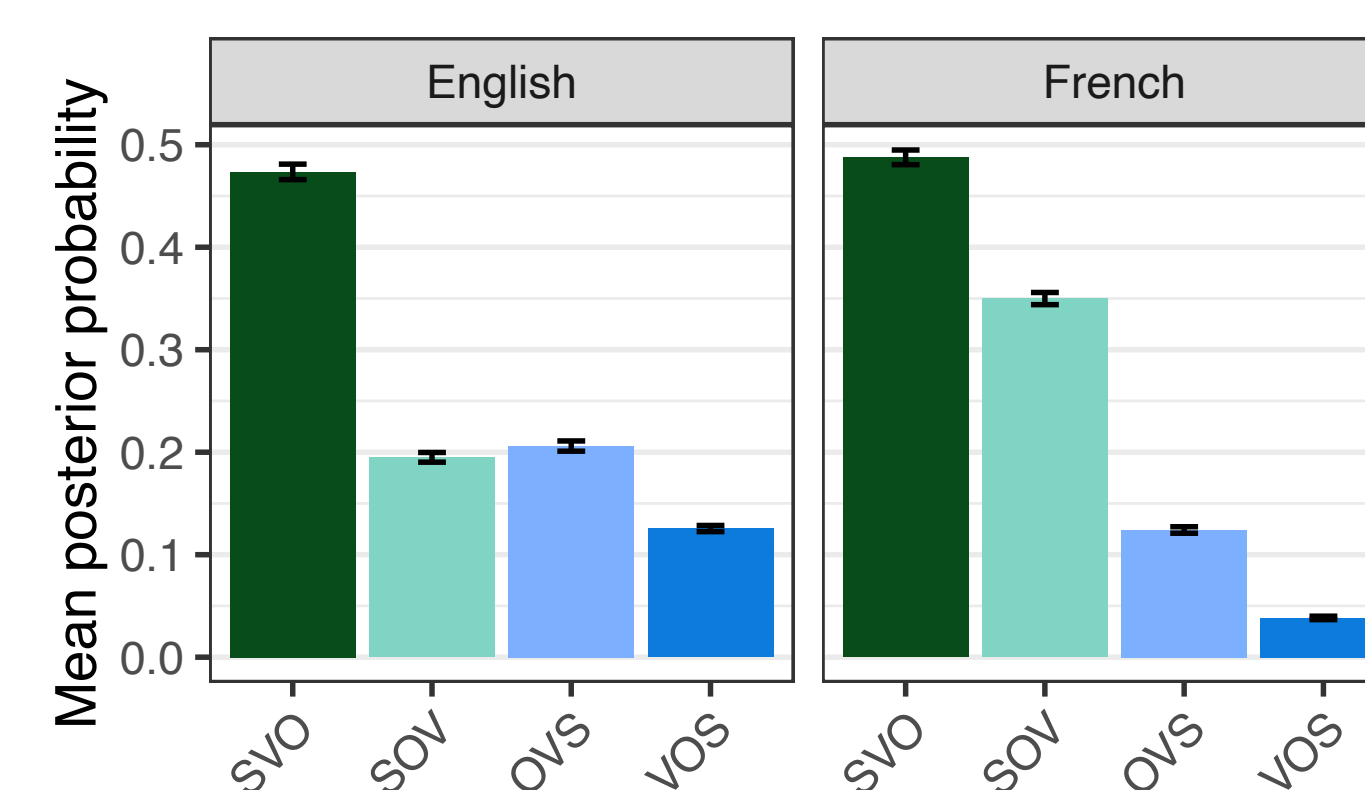


Fig. 2 Posterior distribution over word order grammars

Discussion

We find that input filtering can in principle enable acquisition of basic word order from noisy data

- From imperfectly-identified NP and V distributions alone, our model learns to separate evidence for canonical word order from the distorting effects of "noise" processes
- It does so without knowing ahead of time what noise looks like, or how much there is

Restrictive options in the learner's hypothesis space allow successful filtering

- Each word order grammar allows only a certain combination of rules
- Preference emerges to use these when possible, rather than analyzing everything as noise

- Provides a novel alternative to **regularization** in grammar learning

- Grammar leads you to expect regularities in your data
 ► Filtering allows you to find them

- What do the data from the canonical grammar look like?
- What do the data from noise look like?
- What is the right division into signal vs. noise?

