# [TACL] Surface Statistics of an Unknown Language Indicate How to Parse It
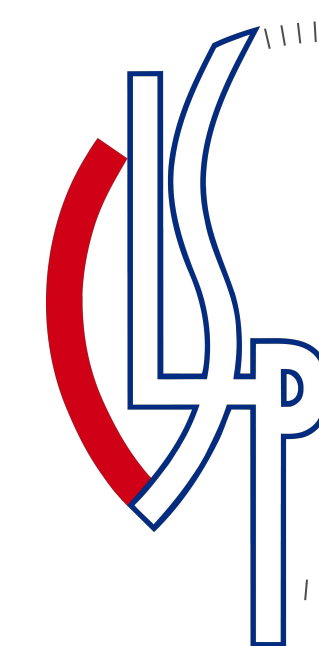
Dingquan Wang and Jason Eisner
{wdd,eisner}@jhu.edu
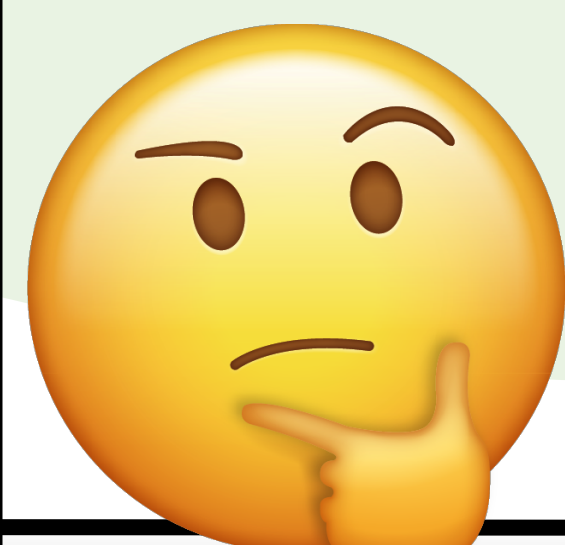
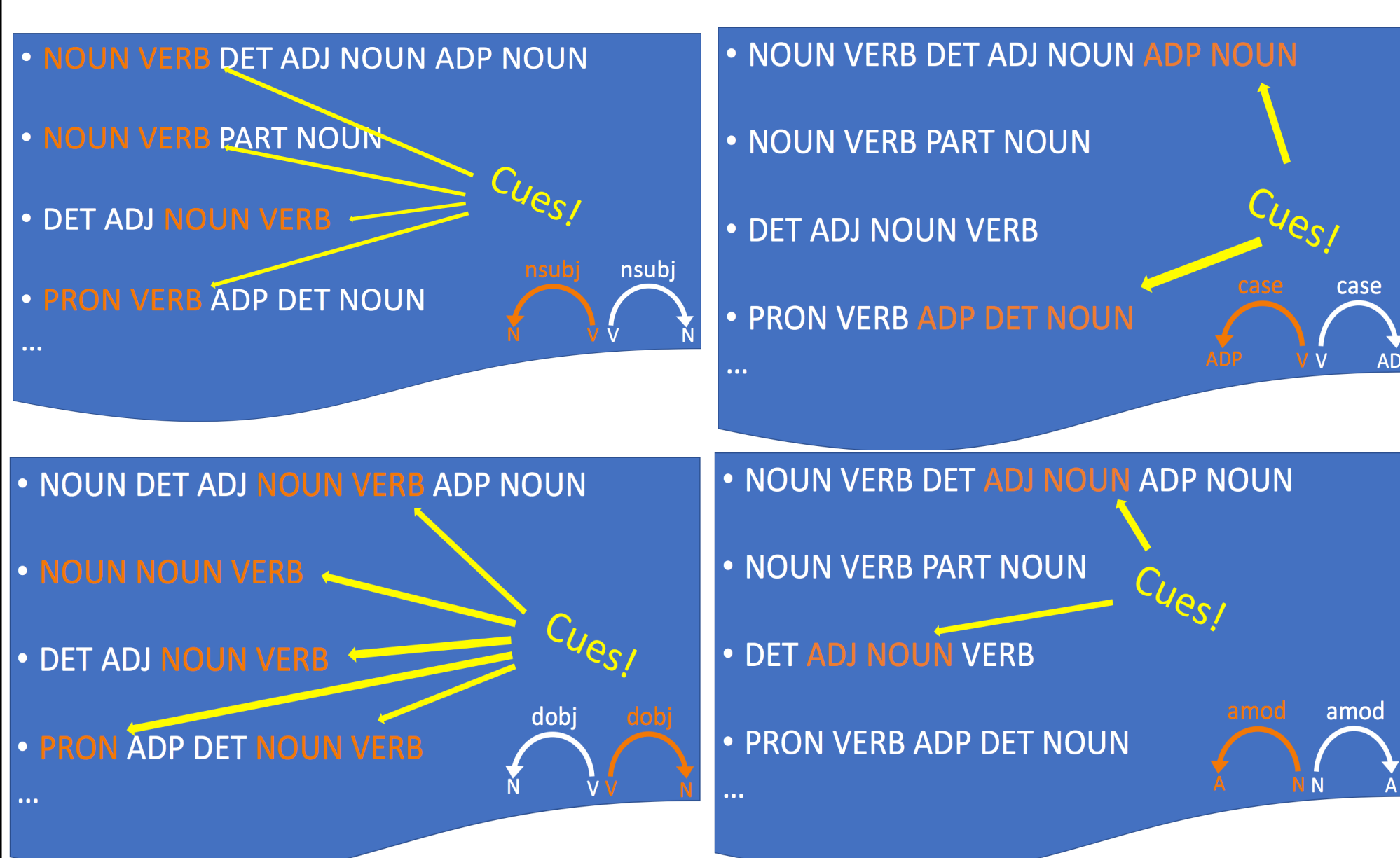You find a **POS-tagged corpus** of text in an *unknown* language.

Can you parse this?
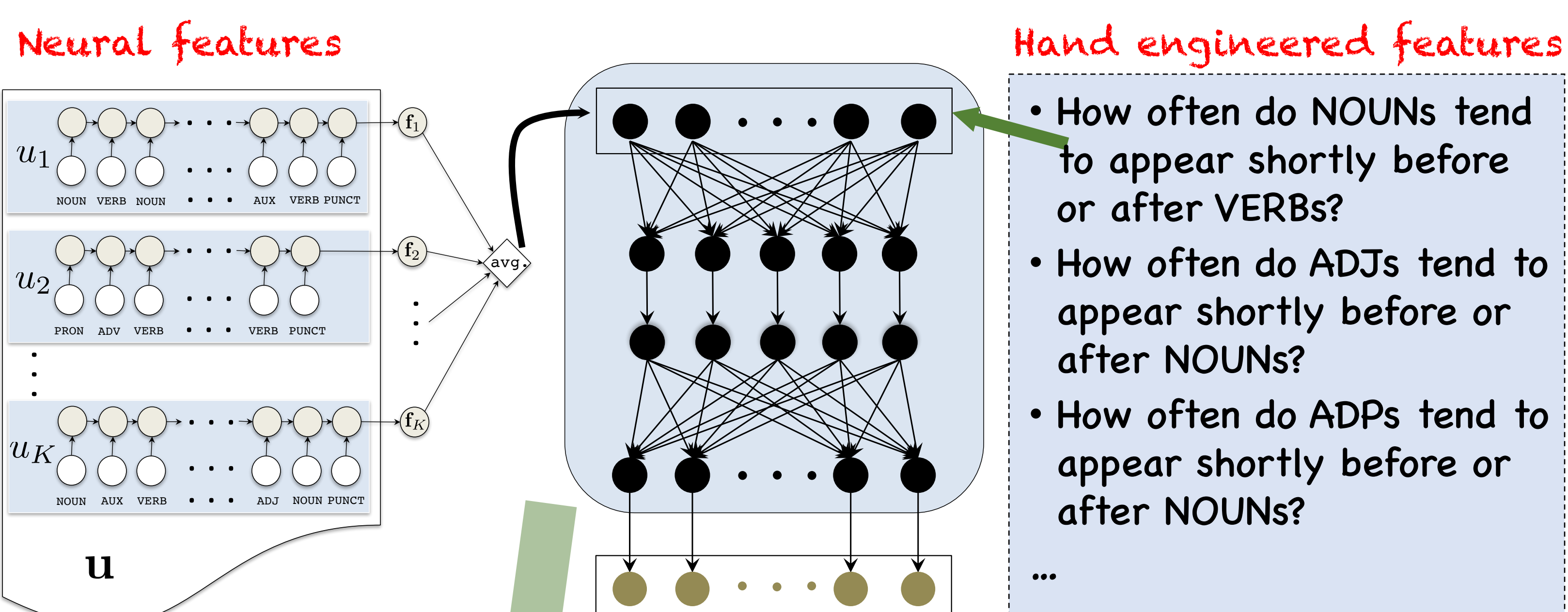
**VERB DET NOUN ADJ DET NOUN**

- Let's extract interesting features of the whole corpus ("surface cues to structure").
- Our universal parser sees these corpus features, along with the input sentence.
- The universal parser is trained end-to-end on diverse languages, with supervision from treebanks.
- Including treebanks for thousands of *synthetic training languages*. This helps.
- Our best method improved UAS and LAS on held-out test languages by an average of **5.6** percentage points over past work.

## Each language is an example!
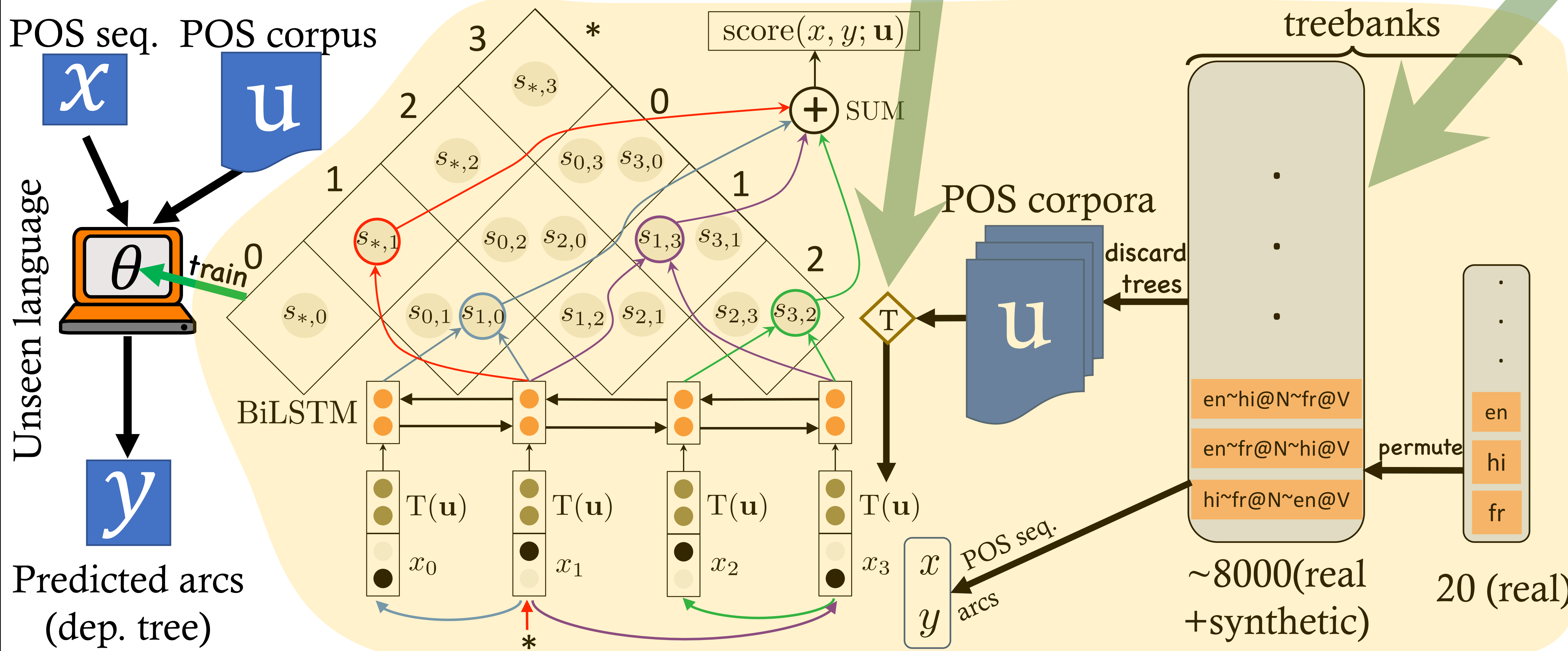
POS corpus → u

True Parser

Language 1 ( [PRON AUX ... / VERB PROPN ...] , Parser1 )

Language 2 ( [VERB PROPN ... / NOUN DET ... / NOUN ADJ ...] , Parser2 )

POS corpus → u

Treebank

Language 1 ( [PRON AUX ... / VERB PROPN ...] , [...] )

Language 2 ( [VERB NOUN... / NOUN DET... / NOUN ADJ...] , [...] )

## Surface Cues to Structure

NOUN VERB DET ADJ NOUN ADP NOUN
NOUN VERB PART NOUN
DET ADJ NOUN VERB
PRON VERB ADP DET NOUN
...
Cues!
nsubj / nsubj

NOUN VERB DET ADJ NOUN ADP NOUN
NOUN VERB PART NOUN
DET ADJ NOUN VERB
PRON VERB ADP DET NOUN
...
Cues!
case / case

NOUN DET ADJ NOUN VERB ADP NOUN
NOUN NOUN VERB
DET ADJ NOUN VERB
PRON ADP DET NOUN VERB
...
Cues!
dobj / dobj

NOUN VERB DET ADJ NOUN ADP NOUN
NOUN VERB PART NOUN
DET ADJ NOUN VERB
PRON VERB ADP DET NOUN
...
Cues!
amod / amod

## The Typology Component

**Neural features**

$u_1$ NOUN VERB NOUN ... AUX VERB PUNCT $f_1$
$u_2$ PRON ADV VERB ... VERB PUNCT $f_2$
$u_K$ NOUN AUX VERB ... ADJ NOUN PUNCT $f_K$
avg
**u**

**Hand engineered features**

- How often do NOUNs tend to appear shortly before or after VERBs?
- How often do ADJs tend to appear shortly before or after NOUNs?
- How often do ADPs tend to appear shortly before or after NOUNs?
- ...

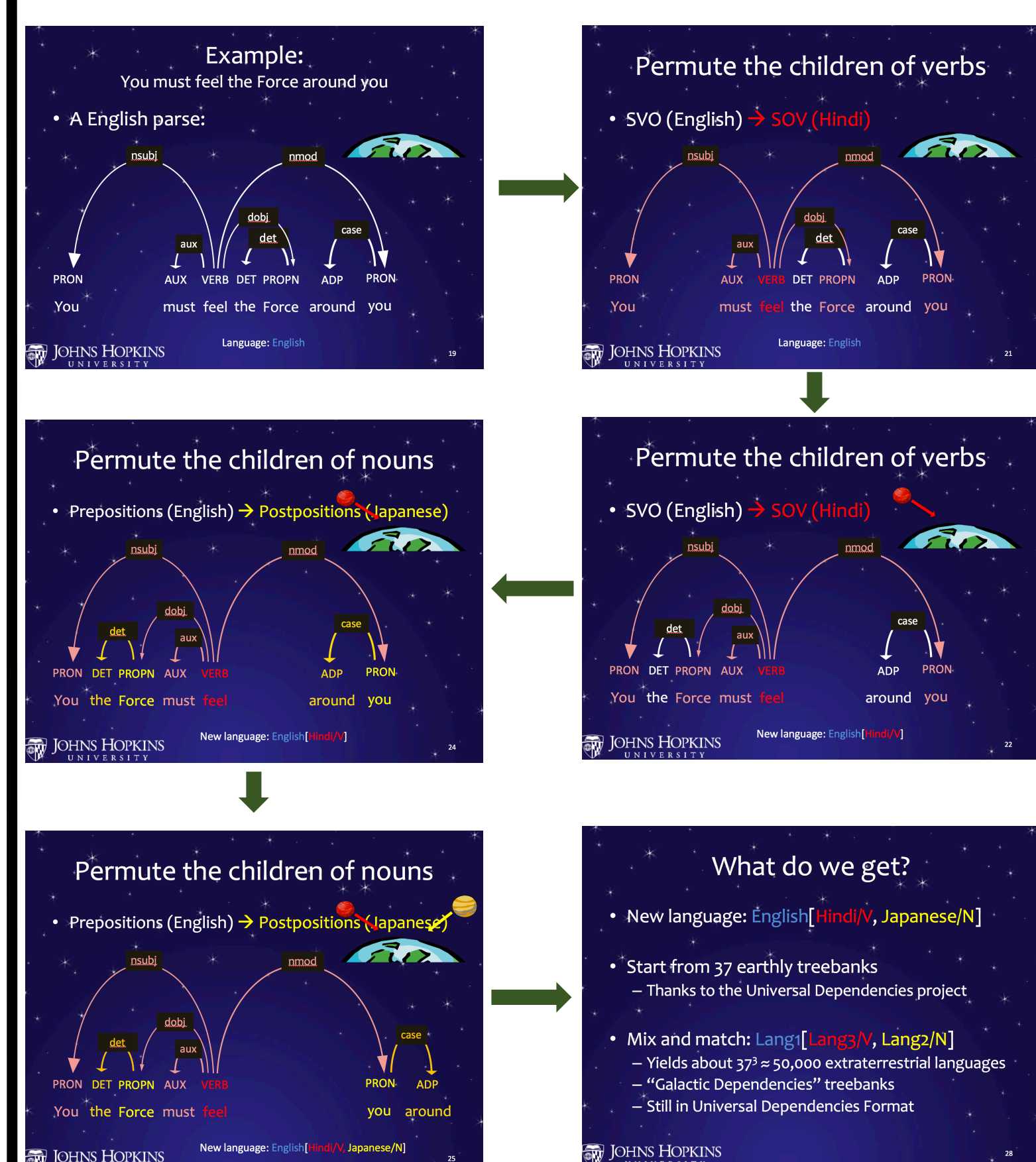## Galactic Treebanks (Wang & Eisner 2016)

- More than 50,000 synthetic languages
  - Resemble real languages, but not found on Earth
- Each has a corpus of dependency parses
  - In the Universal Dependencies format
  - Vertices are words labeled with POS tags
  - Edges are labeled syntactic relationships
- Provide train/dev/test splits, alignments, tools

## What are we building?

POS seq. $x$

POS corpus $u$

$\theta$ — train

Unseen language

Predicted arcs (dep. tree) $y$

## How do we train?

$\text{score}(x, y; \mathbf{u})$

SUM

$s_{*,3}$  3  *
$s_{*,2}$  2   $s_{0,3}$ 0  $s_{3,0}$ 3
$s_{*,1}$  1  $s_{0,2}$ $s_{2,0}$  $s_{1,3}$ $s_{3,1}$ 1
$s_{*,0}$  0  $s_{0,1}$ $s_{1,0}$ $s_{1,2}$ $s_{2,1}$ $s_{2,3}$ $s_{3,2}$ 2

BiLSTM

$T(\mathbf{u})$ $x_0$ | $T(\mathbf{u})$ $x_1$ | $T(\mathbf{u})$ $x_2$ | $T(\mathbf{u})$ $x_3$

T

POS corpora $u$ — discard trees

treebanks

$x$ / $y$ — POS seq. / arcs

~8000 (real + synthetic)

en / hi / fr — permute

20 (real)

en~hi@N~fr@V
en~fr@N~hi@V
hi~fr@N~en@V

## Galactic Dependencies example slides

Example: You must feel the Force around you

Permute the children of verbs — SVO (English) → SVO [Mind]

Permute the children of nouns — Prepositions (English) → Postpositions (Japanese)

Permute the children of verbs — SVO (English) → SOV [Mind]

Permute the children of nouns — Prepositions (English) → Postpositions (Japanese)

What do we get?
- New language: English [Mind], Japanese[N]
- Start from 37 earthly treebanks
  - Thanks to the Universal Dependencies project
- Mix and match: Lang1[Lang1V, Lang2N]
  - Yields about 37² = 50,000 extraterrestrial languages
  - "Galactic Dependencies" treebanks
  - Still in Universal Dependencies Format

## Results (each bar stretches from labeled to unlabeled score)

Attachment Scores

Basque, Croatian, Greek, Hebrew, Hungarian, Indonesian, Irish, Japanese, Slavonic, Persian, Polish, Romanian, Slovenian, Swedish, Tamil, Avg.

Precision \ Recall \ F1 — Avg. proportion (%)

Legend: △ Precision, ● Recall, Baseline, + T(**u**), + T(**u**) + synthetic

nmod, case, punct, det, nsubj, root, amod, dobj, advmod, conj, cc, mark, aux, acl, advcl, compound, cop, nummod, name, xcomp