# Geometry-Aware Supertagging

## with Heterogeneous Dynamic Convolutions

Konstantinos Kogkalidis[1,2]    &    Michael Moortgat[2]

1 Aalto University

2 Utrecht Institute of Linguistics OTS, Utrecht University

LSD, Sept 12, Göteborg

**Utrecht University**

Categorial Grammars 101

**what are they?**

A **family** of syntactic formalisms; each instance consists of:

- ▶ a **lexicon**
  a map assigning *categories* to words: (quasi-)logical formulas (or ADTs)

- ▶ a small set of **inference rules**
  ways to combine and reduce *expressions* based on their categories

Categorial Grammars 101

**Many variations**: TLG, ACG, CCG, …(*CG)

**common points**

▶ **Lexicalized**
   words come packed with their combinatorics

▶ **Formal**
   proximal to logics, type theory & functional programming

▶ **Transparent**
   neat syntax-semantics interface

CATEGORIAL GRAMMARS
○●○

SUPERTAGGING
○○○○

GEOMETRY
○○○○○○

Categorial Grammars 101

**Many variations**: TLG, ACG, CCG, …(*CG)

**divergences**
different background logics $\implies$

- ▶ different linguistic aspects captured
  *e.g. surface order, non-local syntax, dependency relations*

- ▶ different parsing complexity

- ▶ different computational semantics

- ▶ …

CATEGORIAL GRAMMARS
○○●

SUPERTAGGING
○○○○

GEOMETRY
○○○○○○

## Categorial Grammars 101

but! the **parsing pipeline** is always the same
given an input sentence:

1. Assign a category to each word
2. Build the syntactic derivation bottom-up
3. ???
4. Profit

CATEGORIAL GRAMMARS
○○○

SUPERTAGGING
●○○○

GEOMETRY
○○○○○○

Supertagging: the task

For some input sentence $w_1, \ldots w_n$ find the category assignment $c_1, \ldots c_n$ s.t.

$$argmax_{(c_1, \ldots c_n)} \; p(c_1, \ldots c_n \mid w_1, \ldots w_n)^*$$

CATEGORIAL GRAMMARS
OOO

SUPERTAGGING
●OOO

GEOMETRY
OOOOOO

Supertagging: the task

For some input sentence $w_1, \ldots w_n$ find the category assignment $c_1, \ldots c_n$
s.t.

$$argmax_{(c_1, \ldots c_n)} \ p(c_1, \ldots c_n \mid w_1, \ldots w_n)^*$$

*In practice:
build the best statistical model possible given current technology and available
data

## (pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
  *seq2seq (late 10s)*

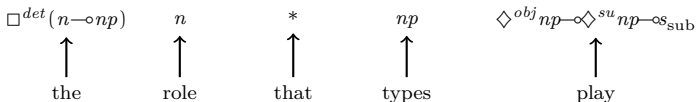| $\Box^{det}(n \multimap np)$ | $n$ | $*$ | $np$ | $\Diamond^{obj} np \multimap \Diamond^{su} np \multimap s_{\text{sub}}$ |
|---|---|---|---|---|
| the | role | that | types | play |

$$* := \Diamond^{body}(\Diamond^{obj1} np \multimap s_{\text{sub}}) \multimap \Box^{mod}(np \multimap np)$$

## (pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
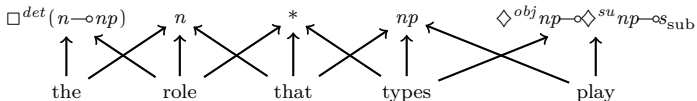  *seq2seq (late 10s)*

| $\Box^{det}(n \multimap np)$ | $n$ | $*$ | $np$ | $\Diamond^{obj} np \multimap \Diamond^{su} np \multimap s_{\text{sub}}$ |
|:---:|:---:|:---:|:---:|:---:|
| ↑ | ↑ | ↑ | ↑ | ↑ |
| the | role | that | types | play |

$* := \Diamond^{body}(\Diamond^{obj1} np \multimap s_{\text{sub}}) \multimap \Box^{mod}(np \multimap np)$

(pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
  *seq2seq (late 10s)*



$* := \diamond^{body} ( \diamond^{obj1} np \multimap s_{\text{sub}} ) \multimap \square^{mod} ( np \multimap np )$

## (pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- ▶ $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- ▶ $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- ▶ $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- ▶ $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
  *seq2seq (late 10s)*

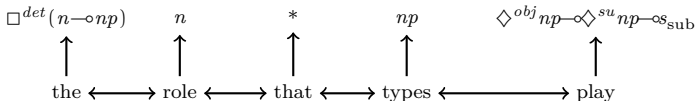$\square^{det}(n \multimap np) \qquad n \qquad * \qquad np \qquad \diamond^{obj} np \multimap \diamond^{su} np \multimap s_{\text{sub}}$

the $\longleftrightarrow$ role $\longleftrightarrow$ that $\longleftrightarrow$ types $\longleftrightarrow$ play

$* := \diamond^{body}(\diamond^{obj1} np \multimap s_{\text{sub}}) \multimap \square^{mod}(np \multimap np)$

(pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
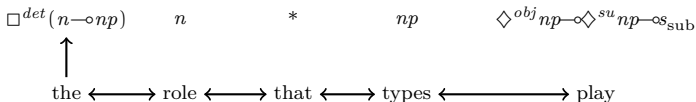  *seq2seq (late 10s)*



$* := \diamondsuit^{body}(\diamondsuit^{obj1}np \multimap s_{sub}) \multimap \square^{mod}(np \multimap np)$

(pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
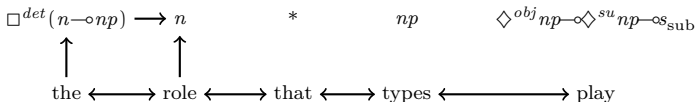  *seq2seq (late 10s)*

$$\square^{det}(n \multimap np) \qquad n \qquad * \qquad np \qquad \diamond^{obj} np \multimap \diamond^{su} np \multimap s_{\mathrm{sub}}$$

$$\uparrow$$

$$\text{the} \longleftrightarrow \text{role} \longleftrightarrow \text{that} \longleftrightarrow \text{types} \longleftrightarrow \text{play}$$

$$* := \diamond^{body}(\diamond^{obj1} np \multimap s_{\mathrm{sub}}) \multimap \square^{mod}(np \multimap np)$$

## (pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
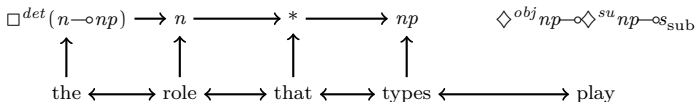  *seq2seq (late 10s)*

$$\Box^{det}(n \multimap np) \longrightarrow n \qquad * \qquad np \qquad \Diamond^{obj} np \multimap \Diamond^{su} np \multimap s_{\text{sub}}$$

$$\uparrow \qquad\qquad \uparrow$$

the $\longleftrightarrow$ role $\longleftrightarrow$ that $\longleftrightarrow$ types $\longleftrightarrow$ play

$$* := \Diamond^{body}(\Diamond^{obj1} np \multimap s_{\text{sub}}) \multimap \Box^{mod}(np \multimap np)$$

(pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
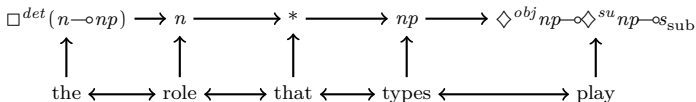  *seq2seq (late 10s)*

$$\square^{det}(n \multimap np) \longrightarrow n \longrightarrow * \qquad np \qquad \Diamond^{obj} np \multimap \Diamond^{su} np \multimap s_{\text{sub}}$$

$$\uparrow \qquad\qquad \uparrow \qquad\qquad \uparrow$$

the $\longleftrightarrow$ role $\longleftrightarrow$ that $\longleftrightarrow$ types $\longleftrightarrow$ play

$$* := \Diamond^{body}(\Diamond^{obj1} np \multimap s_{\text{sub}}) \multimap \square^{mod}(np \multimap np)$$

CATEGORIAL GRAMMARS
000

SUPERTAGGING
0●00

GEOMETRY
000000

(pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
  *seq2seq (late 10s)*

$$\Box^{det}(n \multimap np) \longrightarrow n \longrightarrow * \longrightarrow np \qquad \Diamond^{obj} np \multimap \Diamond^{su} np \multimap s_{\text{sub}}$$

$$\uparrow \qquad\qquad \uparrow \qquad\qquad \uparrow \qquad\qquad \uparrow$$

$$\text{the} \longleftrightarrow \text{role} \longleftrightarrow \text{that} \longleftrightarrow \text{types} \longleftrightarrow \text{play}$$

$$* := \Diamond^{body}(\Diamond^{obj1} np \multimap s_{\text{sub}}) \multimap \Box^{mod}(np \multimap np)$$

(pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
  *seq2seq (late 10s)*

$$\Box^{det}(n{\multimap}np) \longrightarrow n \longrightarrow * \longrightarrow np \longrightarrow \Diamond^{obj}np{\multimap}\Diamond^{su}np{\multimap}s_{\text{sub}}$$

$$\uparrow \qquad \uparrow \qquad \uparrow \qquad \uparrow \qquad \uparrow$$

$$\text{the} \longleftrightarrow \text{role} \longleftrightarrow \text{that} \longleftrightarrow \text{types} \longleftrightarrow \text{play}$$

$$* := \Diamond^{body}(\Diamond^{obj1}np{\multimap}s_{\text{sub}}){\multimap}\Box^{mod}(np{\multimap}np)$$

CATEGORIAL GRAMMARS
○○○

SUPERTAGGING
○●○○

GEOMETRY
○○○○○○

(pre-)history

$p(t_1, \ldots t_n \mid w_1, \ldots w_n) \approx$

- ► $\prod_i^n (t_i \mid w_i)$
  *co-occurrence-based statistical models (90s)*

- ► $\prod_i^n (t_i \mid w_{i-\kappa} \ldots w_{i+\kappa})$
  *window-based n-gram models (00s), FFNs (early 10s)*

- ► $\prod_i^n (t_i \mid w_1, \ldots w_n)$
  *sequence encoders (mid 10s)*

- ► $\prod_i^n (t_i \mid t_1, \ldots t_{i-1}, w_1, \ldots w_n)$
  *seq2seq (late 10s)*

what have we done?

- • more arrows (=more context)
- • auto-regression (price: temporal delay)
- • what about the co-domain?

CATEGORIAL GRAMMARS
○○○

SUPERTAGGING
○○●○

GEOMETRY
○○○○○○

# Intermezzo: the curse(?) of sparsity

The majority of unique categories in common datasets are rare

the "*fix*": ignore rare categories

▶ small penalty in accuracy

▶ less so for coverage..

▶ meta: sparse grammars = bad

the **fix**: decompose categories & build them up during decoding

⚡ unlimited ~~power~~ generalization

▶ meta: sparse grammars = ok

CATEGORIAL GRAMMARS
○○○

SUPERTAGGING
○○●○

GEOMETRY
○○○○○○

# Intermezzo: the curse(?) of sparsity

The majority of unique categories in common datasets are rare

the "*fix*": ignore rare categories
- ▶ small penalty in accuracy
- ▶ less so for coverage..
- ▶ meta: sparse grammars = bad

the **fix**: decompose categories & build them up during decoding
- ⚡ unlimited ~~power~~ generalization
- ▶ meta: sparse grammars = ok

## Intermezzo: the curse(?) of sparsity

The majority of unique categories in common datasets are rare

the "*fix*": ignore rare categories
- ▶ small penalty in accuracy
- ▶ less so for coverage..
- ▶ meta: sparse grammars = bad

the **fix**: decompose categories & build them up during decoding
- ↯ unlimited ~~power~~ generalization
- ▶ meta: sparse grammars = ok

## Modern Times

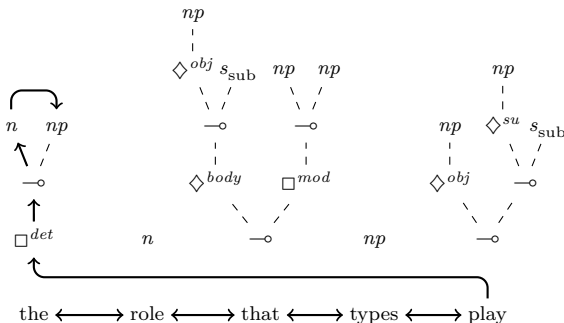$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \text{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

## Modern Times

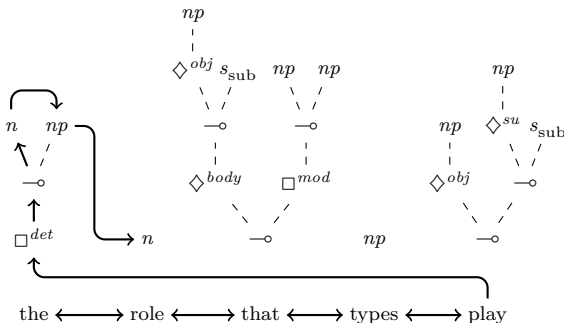$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \text{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*



the ⟷ role ⟷ that ⟷ types ⟷ play

## Modern Times

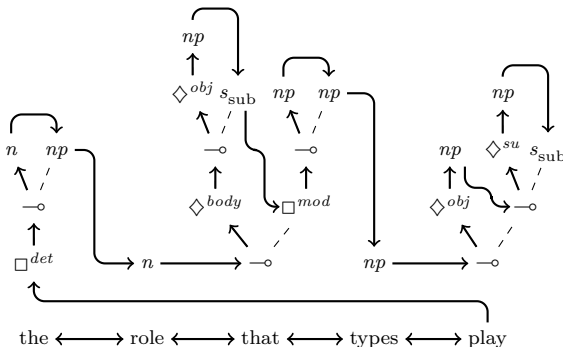$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

► $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

► $\prod_i^m (\sigma_i \mid \mathrm{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

CATEGORIAL GRAMMARS
000

SUPERTAGGING
0000●

GEOMETRY
000000

## Modern Times

$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

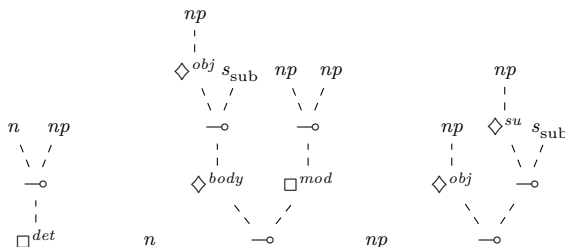- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \mathrm{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

CATEGORIAL GRAMMARS
000

SUPERTAGGING
0000●

GEOMETRY
000000

## Modern Times

$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \mathrm{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

CATEGORIAL GRAMMARS
000

SUPERTAGGING
○○○●

GEOMETRY
000000

## Modern Times

$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \mathrm{anc}(\sigma_i), w_1, \ldots w_n)$
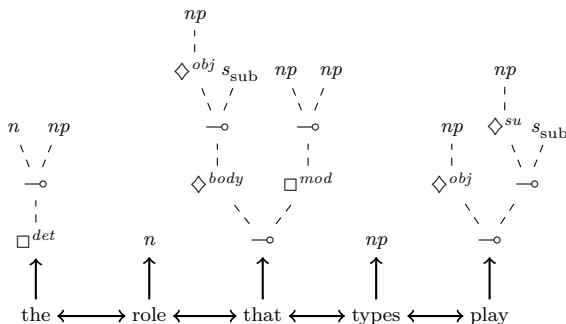  *tree-recursive (Prange et. al 2020)*

## Modern Times

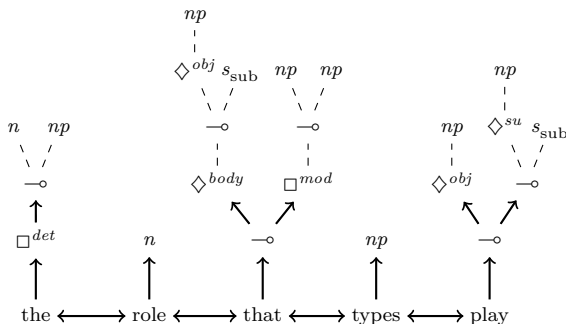$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \text{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

## Modern Times

$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \mathrm{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

## Modern Times

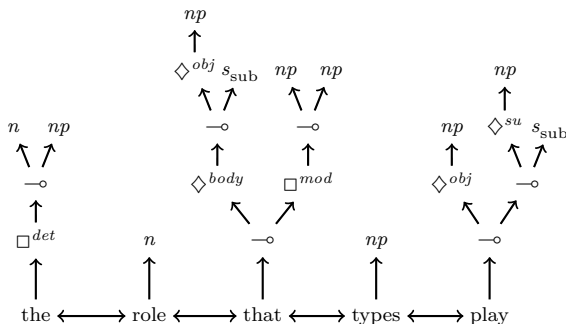$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \mathrm{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

## Modern Times

$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$
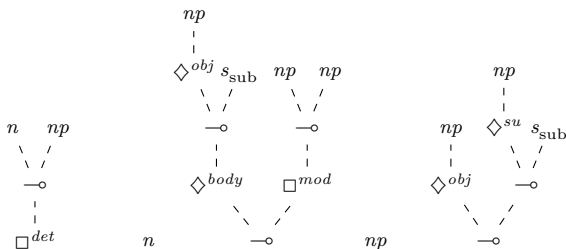
- ▶ $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- ▶ $\prod_i^m (\sigma_i \mid \mathrm{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

## Modern Times

$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m (\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m (\sigma_i \mid \text{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*

CATEGORIAL GRAMMARS
000

SUPERTAGGING
000●

GEOMETRY
000000

## Modern Times

$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx$

- $\prod_i^m(\sigma_i \mid \sigma_1, \ldots \sigma_{i-1}, w_1, \ldots w_n)$
  *sequential constructive (w/ Moortgat & Deoskar, 2019)*

- $\prod_i^m(\sigma_i \mid \operatorname{anc}(\sigma_i), w_1, \ldots w_n)$
  *tree-recursive (Prange et. al 2020)*
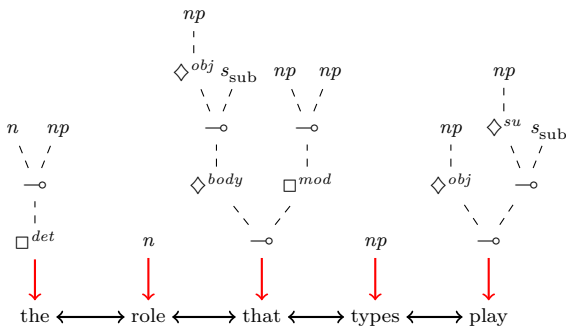
CATEGORIAL GRAMMARS
ooo

SUPERTAGGING
oooo

GEOMETRY
●oooooo

## Post-modernity

neither sequence nor tree but sequence of trees

$$p(\sigma_1, \dots \sigma_m \mid w_1, \dots w_n) \approx \prod_i^m (\sigma_i \mid \sigma_j : \mathrm{depth}(\sigma_j) < \mathrm{depth}(\sigma_i), w_1, \dots w_n)$$



*(encode)*

## Post-modernity

neither sequence nor tree but sequence of trees

$$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx \prod_i^m(\sigma_i \mid \sigma_j : \mathrm{depth}(\sigma_j) < \mathrm{depth}(\sigma_i), w_1, \ldots w_n)$$



*(predict)*

## Post-modernity

neither sequence nor tree but sequence of trees

$$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx \prod_i^m (\sigma_i \mid \sigma_j : \text{depth}(\sigma_j) < \text{depth}(\sigma_i), w_1, \ldots w_n)$$
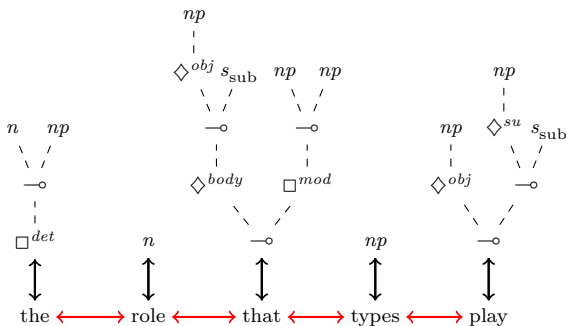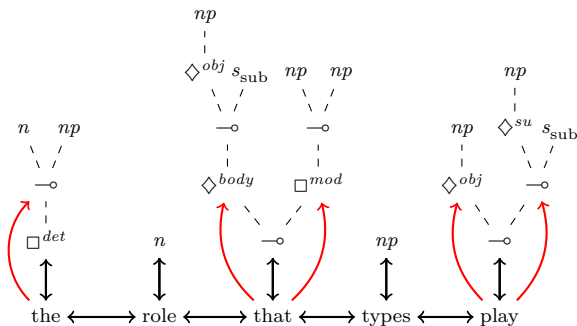


*(feedback)*

## Post-modernity

neither sequence nor tree but sequence of trees

$$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx \prod_i^m (\sigma_i \mid \sigma_j : \mathrm{depth}(\sigma_j) < \mathrm{depth}(\sigma_i), w_1, \ldots w_n)$$



*(contextualize)*

## Post-modernity

neither sequence nor tree but sequence of trees

$$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx \prod_i^m (\sigma_i \mid \sigma_j : \mathrm{depth}(\sigma_j) < \mathrm{depth}(\sigma_i), w_1, \ldots w_n)$$
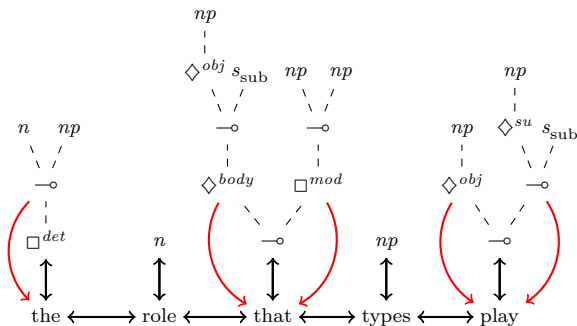


*(predict)*

## Post-modernity

neither sequence nor tree but sequence of trees

$$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx \prod_i^m (\sigma_i \mid \sigma_j : \mathrm{depth}(\sigma_j) < \mathrm{depth}(\sigma_i), w_1, \ldots w_n)$$



*(feedback)*

## Post-modernity

neither sequence nor tree but sequence of trees

$$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx \prod_i^m (\sigma_i \mid \sigma_j : \mathrm{depth}(\sigma_j) < \mathrm{depth}(\sigma_i), w_1, \ldots w_n)$$



*(contextualize)*

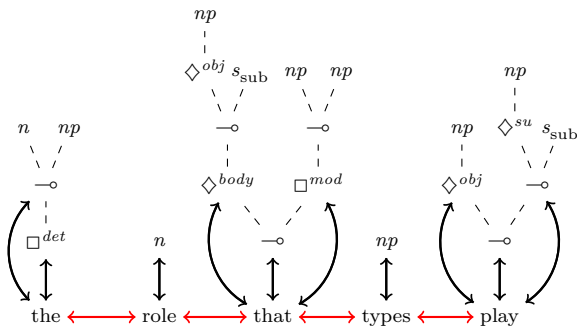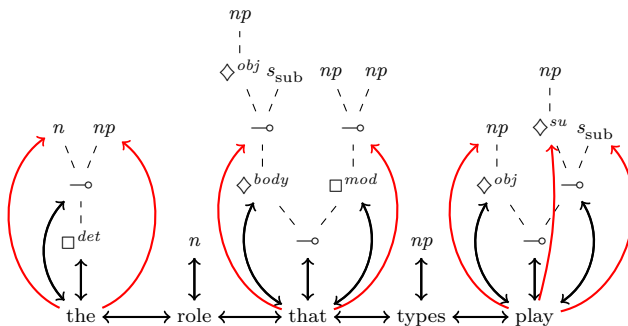## Post-modernity

neither sequence nor tree but sequence of trees

$$p(\sigma_1, \ldots \sigma_m \mid w_1, \ldots w_n) \approx \prod_i^m (\sigma_i \mid \sigma_j : \mathrm{depth}(\sigma_j) < \mathrm{depth}(\sigma_i), w_1, \ldots w_n)$$



*(predict)*

Implementation: dynamic graph convolutions

1 decoding step per tree depth; 3 message-passing rounds per step

▶ *contextualize: states → states*
universal transformer encoder w/ relative weights
(many-to-many, update states with neighborhood context)

▶ *predict: state → nodes*
token classification w/ dynamic tree embeddings
(one-to-many, predict fringe nodes from current state)

▶ *feedback: nodes → state*
heterogeneous graph attention
(many-to-one, update state with last predicted nodes)

CATEGORIAL GRAMMARS
000

SUPERTAGGING
0000

GEOMETRY
000●000

## Table with numbers

| model | accuracy (%) | | | | |
|---|---|---|---|---|---|
| | overall | frequent | uncommon | rare | unseen |
| **CCGbank** (Combinatory Categorial Grammar, en) | | | | | |
| Sequential RNN | 95.10 | 95.48 | 65.76 | 26.02 | 0.00 |
| Tree Recursive | 96.09 | 96.44 | 68.10 | 37.40 | 3.03 |
| Attentive Convolutions | 96.25 | 96.64 | 71.04 | – | – |
| *this work* | 96.29 | 96.61 | 72.06 | 34.45 | 4.55 |
| **CCGrebank** (ditto, improved version) | | | | | |
| Sequential RNN | 94.44 | 94.93 | 66.90 | 27.41 | 1.23 |
| Tree Recursive | 94.70 | 95.11 | 68.86 | 36.76 | 4.94 |
| *this work* | 95.07 | 95.45 | 71.40 | 37.19 | 3.70 |
| **TLGBank** (Lambek calculus & control modalities, fr) | | | | | |
| ELMo LSTM | 93.20 | 95.10 | 75.19 | 25.85 | – |
| *this work* | 95.93 | 96.40 | 81.48 | 55.37 | 7.26 |
| **Æthel** (van Benthem calculus & dependency modalities, nl) | | | | | |
| Sequential Transformer | 83.67 | 84.55 | 64.70 | 50.58 | 24.55 |
| *this work* | 93.67 | 94.72 | 73.45 | 53.83 | 15.78 |

## Color coded summary

| *decoder* | seq2seq[*t*] | seq2seq[σ] | tree | dynamic graph |
|---|---|---|---|---|
| *codomain* | fixed | open | constrained | constrained |
| *context* | left | preorder (global) | ancestors (local) | levels (global) |
| *complexity* | # words | # symbols | tree depth | tree depth |
| *treeness* | ignored | implicit | explicit | explicit |
| *sequencess* | explicit | misaligned | ignored | explicit |
| *search?* | ✓ | ✓ | ? | ? |

**legend**

▶ green = good

▶ yellow = meh

▶ red = bad

Take home messages

use hammers for nails only

sparsity.. a *friend*?

▶ more rare cats $\implies$ better acquisition of rare cats
▶ cascading effect on performance

thanks!