

Context-free grammars & finite-state automata over categories*

Noam Zeilberger

joint work with Paul-André Mellies

PSSL 109 @ Leiden
15-17 November 2024

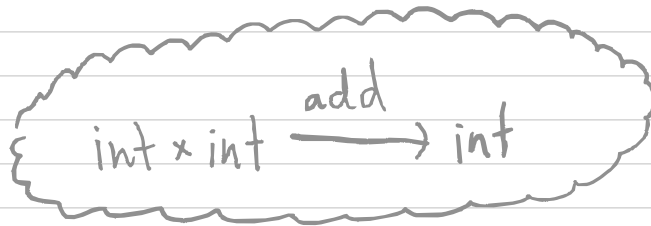
* Paper: The categorical contours of the Chomsky-Schützenberger Representation Theorem

arXiv:2405.14703
(new version 15 Nov!)

Typing as a lifting problem

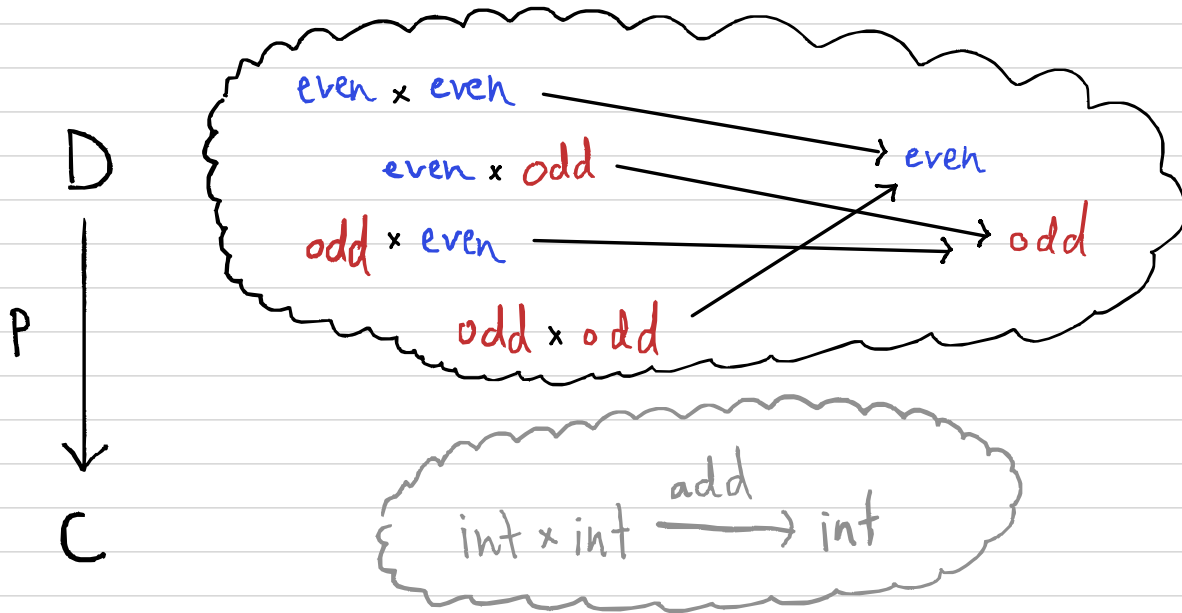
Idea: model type systems fibrationally as functors that "forget" typing information.

C



Typing as a lifting problem

Idea: model type systems fibrationally as functors that "forget" typing information.



See "Functors are Type Refinement Systems" (POPL 2015) and other papers in series w/ PAM.

~~Typing~~^{Parsing} as a lifting problem

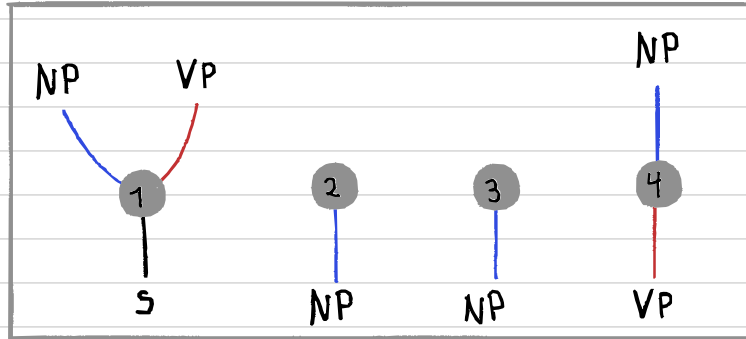
$S \rightarrow NP VP$

$NP \rightarrow \text{Noam}$

$NP \rightarrow \text{PSSL}$

$VP \rightarrow \text{likes } NP$

~~Typing~~^{Parsing} as a lifting problem



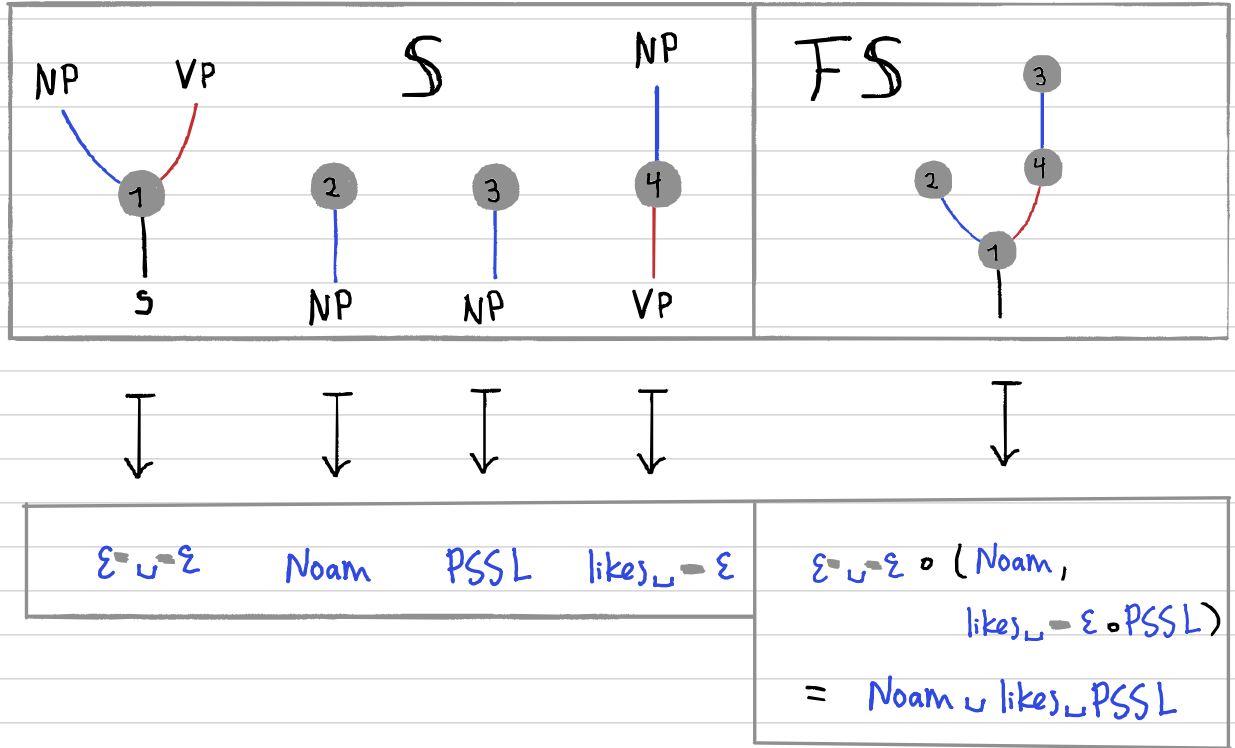
- 1: $S \rightarrow NP \ VP$
- 2: $NP \rightarrow \text{Noam}$
- 3: $NP \rightarrow \text{PSSL}$
- 4: $VP \rightarrow \text{likes } NP$

$\varepsilon \rightarrow \varepsilon$ Noam PSSL likes_⊥ - ε

~~Typing~~^{Parsing} as a lifting problem

- 1: $S \rightarrow NP VP$
- 2: $NP \rightarrow \text{Noam}$
- 3: $NP \rightarrow \text{PSSL}$
- 4: $VP \rightarrow \text{likes } NP$

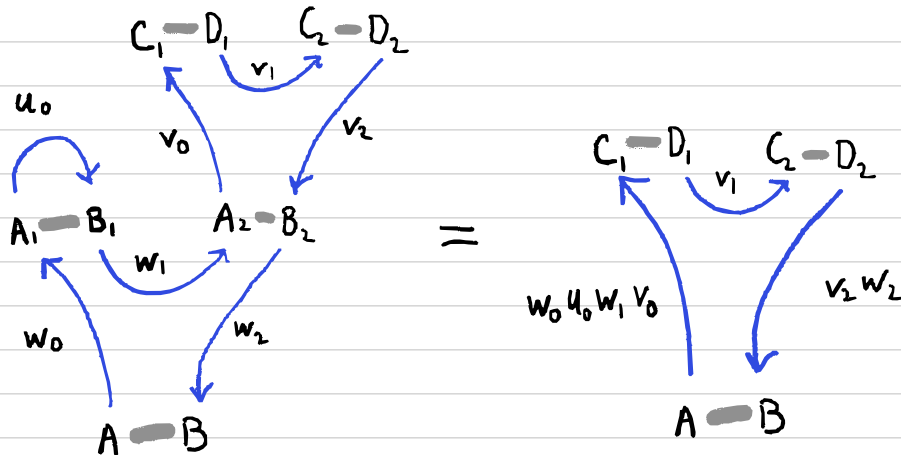
FS
 \downarrow
P
 \downarrow
W Σ



The spliced arrows construction

Given a category \mathcal{C} , the **operad of spliced arrows** WC has:

- Objects given by pairs (A, B) of objects $A, B \in \mathcal{C}$
- n -ary operations $f: (A_1, B_1), \dots, (A_n, B_n) \rightarrow (A, B)$ given by sequences $f = w_0 - \dots - w_n$ of $n+1$ arrows $w_i: B_i \rightarrow A_{i+1} \in \mathcal{C}$
(under convention $B_0 = A, A_{n+1} = B$)
- Composition performed by "splicing into the gaps"...



The spliced arrows construction

Splicing extends to a functor

$$\text{Cat} \xrightarrow{W} \text{Oper}$$

since any functor of categories $F: C \rightarrow D$ induces a functor of operads $WF: WC \rightarrow WD$ by the mappings

$$(A, B) \mapsto (FA, FB)$$

$$w_0 - \dots - w_n \mapsto Fw_0 - \dots - Fw_n$$

Context-free grammar over a category

Definition. A CFG over a category \mathcal{C} is a pair of a pointed finite species $(\mathcal{S}, S \in \mathcal{S})$ and a functor $p: F\mathcal{S} \rightarrow \mathcal{WC}$.

The context-free language of arrows generated by a CFG $G = (\mathcal{S}, S, p)$ is the set $\mathcal{L}_G = \{p(\alpha) \mid \alpha: S\} \subseteq \mathcal{C}(A, B)$ where $(A, B) = p(S)$.

Context-free grammar over a category

Definition. A CFG over a category \mathcal{C} is a pair of a pointed finite species $(\mathcal{S}, S \in \mathcal{S})$ and a functor $p: F\mathcal{S} \rightarrow \mathcal{WC}$.

The context-free language of arrows generated by a CFG $G = (\mathcal{S}, S, p)$ is the set $\mathcal{L}_G = \{p(\alpha) \mid \alpha: S\} \subseteq \mathcal{C}(A, B)$ where $(A, B) = p(S)$.

Proposition. $L \subseteq \Sigma^*$ is context-free in the classical sense iff $\overset{\alpha \in \Sigma}{\curvearrowright}$ it is the language of arrows of a CFG over $F\mathcal{B}_\Sigma$, where $\mathcal{B}_\Sigma = \overset{\alpha \in \Sigma}{\curvearrowright} *$.

Context-free grammar over a category

Definition. A CFG over a category \mathcal{C} is a pair of a pointed finite species $(\mathcal{S}, S \in \mathcal{S})$ and a functor $p: F\mathcal{S} \rightarrow \mathcal{WC}$. The context-free language of arrows generated by a CFG $G = (\mathcal{S}, S), p$ is the set $\mathcal{L}_G = \{p(\alpha) \mid \alpha: S\} \subseteq \mathcal{C}(A, B)$ where $(A, B) = p(S)$.

Proposition. $L \subseteq \Sigma^*$ is context-free in the classical sense iff it is the language of arrows of a CFG over FB_Σ , where $B_\Sigma = \begin{array}{c} a \in \Sigma \\ \curvearrowright \\ * \end{array}$.

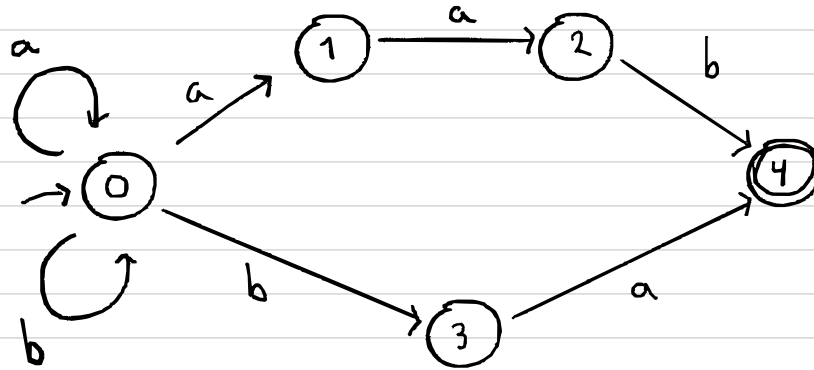
Example #2: Let $B_\Sigma^{1\$} = \begin{array}{c} a \in \Sigma \\ \curvearrowright \\ \perp \xrightarrow{\$} * \xrightarrow{\$} T \end{array}$. A CFG over $FB_\Sigma^{1\$}$ can have productions that are only applicable at beginning/end of string.

$$S \rightarrow E\$ \quad (\text{cf. Knuth 1965})$$

Some closure properties of CFLs

- If $L_1 \subseteq C(A, B)$ and $L_2 \subseteq C(A, B)$ are context-free
union then so is $L_1 \cup L_2 \subseteq C(A, B)$
- If $L_1 \subseteq C(A_1, B_1), \dots, L_n \subseteq C(A_n, B_n)$ are context-free
concatenation then so is $w_0 L_1 w_1 \dots L_n w_n \subseteq C(A, B)$
for any $w_0 \dots w_n : (A_1, B_1), \dots, (A_n, B_n) \rightarrow (A, B) \in \mathcal{WC}$
- If $L \subseteq C(A, B)$ is context-free
homomorphic image then so is $F(L) \subseteq D(FA, FB)$
for any functor $F: C \rightarrow D$

~~Typing~~^{Recognition} as a lifting problem

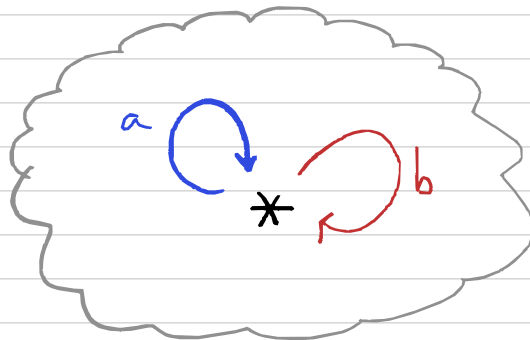
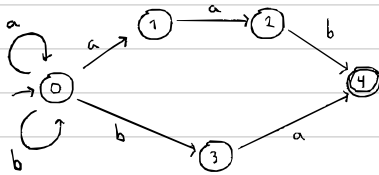
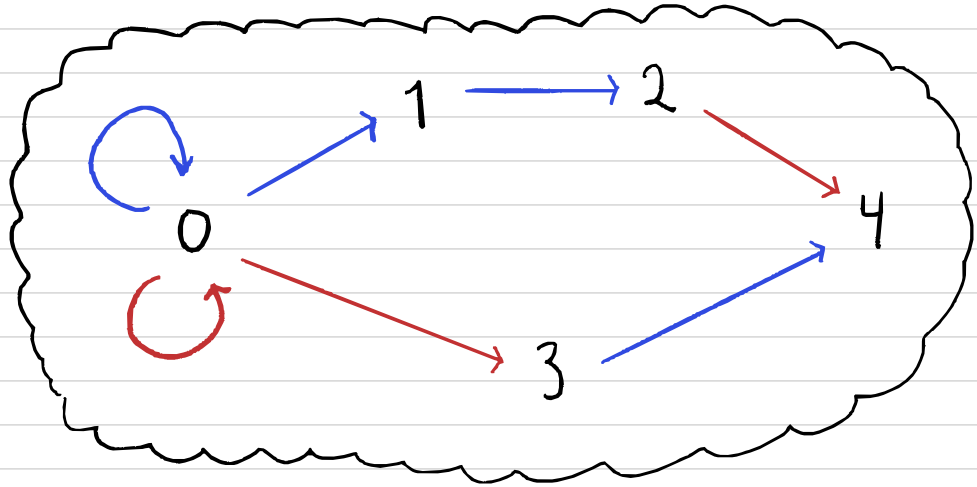


~~Typing~~^{Recognition} as a lifting problem

FQ



FB_Σ



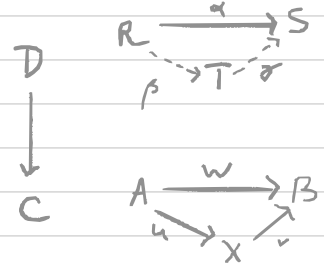
ULF and finitary functors

Let $p: D \rightarrow C$ be a functor of categories.

- p has **unique liftings of factorizations** (ULF aka "discrete Conduché")

if for all arrows $\alpha \in D$ s.t. $p(\alpha) = uv$, $\exists! \beta, \gamma$.

such that $\alpha = \beta\gamma$ and $p(\beta) = u$ and $p(\gamma) = v$.



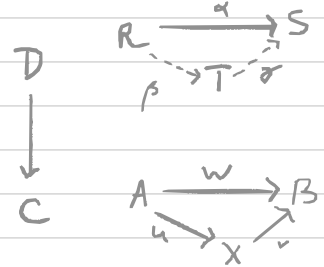
ULF and finitary functors

Let $p: D \rightarrow C$ be a functor of categories.

- p has **unique liftings of factorizations** (ULF aka "discrete Conduché")

if for all arrows $\alpha \in D$ s.t. $p(\alpha) = uv$, $\exists! \beta, \gamma$.

such that $\alpha = \beta\gamma$ and $p(\beta) = u$ and $p(\gamma) = v$.



- p is **finitary** if the fibers $p^{-1}(A)$ and $p^{-1}(w)$ of every object $A \in C$ and every arrow $w: A \rightarrow B$ are finite.

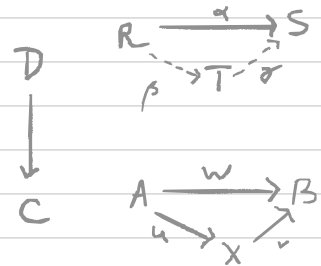
ULF and finitary functors

Let $p: \mathcal{D} \rightarrow \mathcal{C}$ be a functor of categories.

- p has **unique liftings of factorizations** (ULF aka "discrete Conduché")

if for all arrows $\alpha \in \mathcal{D}$ s.t. $p(\alpha) = uv$, $\exists! \beta, \gamma$.

such that $\alpha = \beta\gamma$ and $p(\beta) = u$ and $p(\gamma) = v$.



- p is **finitary** if the fibers $p^{-1}(A)$ and $p^{-1}(w)$ of every object $A \in \mathcal{C}$ and every arrow $w: A \rightarrow B$ are finite.

Proposition. Let $F: \mathcal{C} \rightarrow \text{Span}(\text{Set})$ be the lax functor canonically representing the functor $p: \mathcal{D} \rightarrow \mathcal{C}$. (As a "displayed category".)

- p is ULF iff F is a pseudofunctor
- p is finitary iff F factors via $\text{Span}(\text{FinSet})$

ULF and finitary functors

Proposition (Street 1996, cf. Guetta 2020). Let $p: \mathcal{D} \rightarrow \mathcal{C}$ be a functor into a category $\mathcal{C} \simeq \mathbf{FG}$ freely generated by some graph G . Then p is ULF iff $\mathcal{D} \simeq \mathbf{FH}$ and $p = F\phi$ for some graph H and homomorphism $\phi: H \rightarrow G$.

Proposition. Let $\phi: H \rightarrow G$ be a homomorphism into a finite graph G . Then ϕ is finitary iff H is finite.

Corollary. A functor $p: \mathcal{Q} \rightarrow \mathbf{FB}_{\Sigma}$ represents the underlying bare automaton of an NFA over Σ iff p is ULF and finitary.

Finite-state automaton over a category

Definition. A **NFA over a category C** is a pair of a bipointed category $(Q, q_0 \in Q, q_f \in Q)$ and a finitary ULF functor $p : Q \rightarrow C$. The **regular language of arrows** recognized by an NFA $M = ((Q, q_0, q_f), p)$ is the set
$$\mathcal{L}_M = \{ p(\alpha) \mid \alpha : q_0 \rightarrow q_f \} \subseteq C(A, B).$$
 where $(A, B) = (p(q_0), p(q_f))$.

Finite-state automaton over a category

Definition. A **NFA over a category C** is a pair of a bipoinded category $(Q, q_0 \in Q, q_f \in Q)$ and a finitary ULF functor $p : Q \rightarrow C$. The **regular language of arrows** recognized by an NFA $M = ((Q, q_0, q_f), p)$ is the set $\mathcal{L}_M = \{ p(\alpha) \mid \alpha : q_0 \rightarrow q_f \} \subseteq C(A, B)$.
where $(A, B) = (p(q_0), p(q_f))$

Proposition. $L \subseteq \Sigma^*$ is regular in the classical sense iff it is the language of arrows of a NFA over $FB_{\Sigma}^{\wedge \$}$.

Finite-state automaton over a category

Definition. A NFA over a category C is a pair of a bipoinded category $(Q, q_0 \in Q, q_f \in Q)$ and a finitary ULF functor $p : Q \rightarrow C$. The regular language of arrows recognized by an NFA $M = ((Q, q_0, q_f), p)$ is the set $\mathcal{L}_M = \{ p(\alpha) \mid \alpha : q_0 \rightarrow q_f \} \subseteq C(A, B)$.
where $(A, B) = (p(q_0), p(q_f))$

Proposition. $L \subseteq \Sigma^*$ is regular in the classical sense iff it is the language of arrows of a NFA over $FB_\Sigma^{\uparrow \$}$.

A NFA M is ^(total) deterministic iff p is a discrete opfibration.

It is codeterministic iff p is a discrete fibration.

Some examples of categorical NFA

Product automaton $M \times M' :=$

$Q \times Q'$
 $\downarrow p \times p'$
 $C \times C'$

w/ initial state (q_0, q'_0)
final state (q_f, q'_f)

$\mathcal{L}_{M \times M'} = \mathcal{L}_M \times \mathcal{L}_{M'}$

Some examples of categorical NFA

Product automaton $M \times M' :=$

$$\begin{array}{c} Q \times Q' \\ \downarrow p \times p' \\ C \times C' \end{array} \quad \begin{array}{l} \text{w/ initial state } (q_0, q'_0) \\ \text{final state } (q_f, q'_f) \end{array}$$

$\mathcal{L}_{M \times M'} = \mathcal{L}_M \times \mathcal{L}_{M'}$

Total automaton $M_{C(A,B)} :=$

$$\begin{array}{c} C \\ \downarrow \text{id} \\ C \end{array} \quad \begin{array}{l} \text{w/ initial state } A \\ \text{final state } B \end{array}$$

$\mathcal{L}_{M_{C(A,B)}} = C(A,B)$

Some examples of categorical NFA

Product automaton $M \times M' :=$

$$\begin{array}{c} Q \times Q' \\ \downarrow p \times p' \\ C \times C' \end{array} \quad \begin{array}{l} \text{w/ initial state } (q_0, q'_0) \\ \text{final state } (q_f, q'_f) \end{array}$$

$\mathcal{L}_{M \times M'} = \mathcal{L}_M \times \mathcal{L}_{M'}$

Total automaton $M_{C(A,B)} :=$

$$\begin{array}{c} C \\ \downarrow \text{id} \\ C \end{array} \quad \begin{array}{l} \text{w/ initial state } A \\ \text{final state } B \end{array}$$

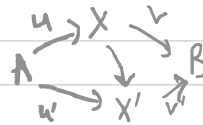
$$\mathcal{L}_{M_{C(A,B)}} = C(A,B)$$

Singleton automaton $M_w :=$

$$\begin{array}{c} \text{Fact}_w \\ \downarrow \tau \\ C \end{array} \quad \begin{array}{l} \text{w/ initial state } (id_A, w) \\ \text{final state } (w, id_B) \end{array}$$

$$\mathcal{L}_{M_w} = \{w\}$$

Requirement: C has finitary factorizations



An aside on ε -transitions

Naturally modelled as arrows $\alpha: q \rightarrow q'$ such that $p(\alpha) = \text{id}$.

... but ULF implies no such arrows! ($p(\text{id}_q \alpha) = p(\alpha \text{id}_{q'}) = \text{id}$)

So to model ε -transitions it seems we need to weaken ULF.

... but the general Conduché property seems too weak.
(wrong version of " ε -removal")

Automata over operads

Definition. NFA over an operad \mathcal{O} is a pair $M = (\dot{Q}, p)$ of a pointed operad $\dot{Q} = (Q, q_r \in Q)$ + a finitary VLF functor $p: Q \rightarrow \mathcal{O}$, recognizing a regular language of constants $\mathcal{L}_M = \{p(a) \mid a: q_r\} \subseteq \mathcal{O}(A)$ where $p(q_r) = A$

Automata over operads

Definition. NFA over an operad \mathcal{O} is a pair $M = (\dot{Q}, p)$ of a pointed operad $\dot{Q} = (Q, q_r \in Q)$ + a finitary VLF functor $p: Q \rightarrow \mathcal{O}$, recognizing a regular language of constants $L_M = \{p(a) \mid a: q_r\} \subseteq \mathcal{O}(A)$ where $p(q_r) = A$

Regular tree language \Leftrightarrow recognized by NFA over free operad.

Automata over operads

Definition. NFA over an operad \mathcal{O} is a pair $M = (\dot{Q}, p)$ of a pointed operad $\dot{Q} = (Q, q_r \in Q)$ + a finitary VLF functor $p: Q \rightarrow \mathcal{O}$, recognizing a regular language of constants $\mathcal{L}_M = \{p(a) \mid a: q_r\} \subseteq \mathcal{O}(A)$ where $p(q_r) = A$

Regular tree language \Leftrightarrow recognized by NFA over free operad.

Proposition. If $p: D \rightarrow C$ is VLF (resp. finitary) then so is $\mathcal{W}p: \mathcal{W}D \rightarrow \mathcal{W}C$

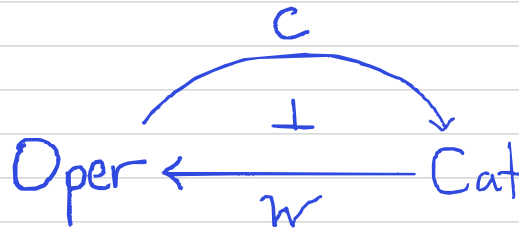
Hence any categorical NFA $M = ((Q, q_o, q_f), p)$ induces an operadic NFA $\mathcal{W}M = ((\mathcal{W}Q, (q_o, q_f)), \mathcal{W}p)$ with

$$\mathcal{L}_{\mathcal{W}M} = \mathcal{L}_M.$$

The Chomsky-Schützenberger Rep Thm (1963)

"A language is context-free iff it is a homomorphic image of a Dyck language with a regular language."

Key observation: the functor \mathcal{W} has a left adjoint!



See our paper for the **contour category** construction, and for the proof of a generalized C-S rep theorem.

One ingredient: closure of CFLs under intersection w/ RLs...

Pulling back a CFG along an NFA

Lemma. The pullback of a functor of operads $p: FS \rightarrow \mathcal{O}$ along a ULF functor of operads $p_Q: \mathcal{Q} \rightarrow \mathcal{O}$ is obtained from a pullback of $\phi: S \rightarrow \mathcal{O}$ along p_Q in species.

$$\begin{array}{ccc}
 S' & \xrightarrow{\gamma'} & S \\
 \downarrow \lrcorner & & \downarrow \phi \\
 \phi' = p_Q^* \phi & & \\
 \downarrow & & \downarrow \\
 \mathcal{Q} & \xrightarrow{p_Q} & \mathcal{O}
 \end{array}
 \quad \Bigg| \quad
 \begin{array}{ccc}
 FS' & \xrightarrow{F\gamma'} & FS \\
 \downarrow \lrcorner & & \downarrow p \\
 p' = p_Q^* p & & \\
 \downarrow & & \downarrow \\
 \mathcal{Q} & \xrightarrow{p_Q} & \mathcal{O}
 \end{array}
 \begin{array}{c}
 \text{Spec} \\
 \text{Oper}
 \end{array}$$

Moreover, if S is finite and p_Q is finitary then S' is finite.

Pulling back a CFG along an NFA

Let $G = ((S, S), p_G)$ be a CFG and $M = ((Q, q_0, q_f), p_M)$ a NFA over the same category w/ $p_G(S) = (p_M(q_0), p_M(q_f))$. Take the pullback:

$$\begin{array}{ccc}
 FS' & \xrightarrow{F\gamma'} & FS \\
 p_G' = Wp_M^* \downarrow p_G & \lrcorner & \downarrow p_G \\
 WQ & \xrightarrow{Wp_M} & WC
 \end{array}
 \qquad
 \begin{array}{c}
 Q \\
 \downarrow p_M \\
 C
 \end{array}$$

Then $M^*G := ((S', (S, (q_0, q_f))), p_G')$ is a CFG generating

$$\mathcal{L}_{M^*G} = p_M^{-1}(\mathcal{L}_G) \cap Q(q_0, q_f).$$

CFL of runs
of the NFA M !

Corollary: CFLs closed under intersection with RLs.

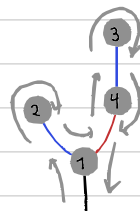
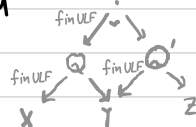
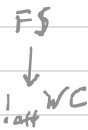
Conclusion

See paper (arXiv:2405.14703) for more on:

- Translations between CFGs
- Generalized CFGs over operads
- More properties of finitary ULF functors, and closure properties of regular languages
- The contour category construction and the universal CFG of a pointed finite species



$$a^n \# b^n \# c^n \in$$



Long term goals:

- Categorify more of automata theory & parsing theory
- Transfer knowledge back to type theory and category theory?