# The ABC's of RASP (mostly ~~B's and~~ C's)
**due to time**
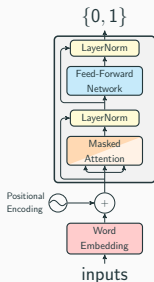
Andy Yang (Univ. of Notre Dame, USA)

11 Oct 2024

# Background

$$(\forall i)(Q_a)$$
$$(\forall i)(\forall j)(i < j \rightarrow Q_a(i) \wedge Q_b(j))$$
etc.

What languages are recognized
by transformer encoders?

What languages are
defined by logical formulas?

For a survey of papers in this area: Strobl et al. [2023], "Transformers as
Recognizers of Formal Languages: A Survey on Expressivity"

## Formal Languages



**Figure 1:** Some complexity classes defined by circuit families, compared with the perhaps more familiar Chomsky hierarchy. See Strobl [2023]

# RASP

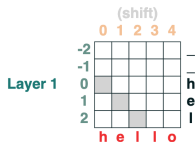**Symbolic Representation of Transformer Computation**

- Matrix multiplications are too confusing for me!
- Can we relate transformers to other formal models?
- How can we prove the expressivity of transformers?

Challenge 2: Shift

Shift all of the tokens in a sequence to the right by i positions. (Here we introduce an optional parameter in the aggregation: the default value to be used when no input positions are selected. If not defined, this value is 0.)

```
def shift(i=1, default="_", seq=tokens):
    x = (key(indices) == query(indices-i)).value(seq, default)
    return x.name("shift")
shift(2)
```



**Figure 2:** Weiss et al. [2021] and https://srush.github.io/raspy/

4

(subset of) RASP $\xrightarrow{\text{[Lindner et al., 2024]}}$ (subset of) transformers

(subset of) RASP $\xleftarrow{\text{[Friedman et al., 2024]}}$ (subset of) transformers

(subset of) RASP $\xrightarrow{\text{[Lindner et al., 2024]}}$ (subset of) transformers

(subset of) RASP $\xleftarrow{\text{[Friedman et al., 2024]}}$ (subset of) transformers

(subset of) RASP $\xrightarrow{\text{[Yang and Chiang, 2024]}}$ transformers

# C-RASP

$$\# \, [j \leq i] \, P(j)$$

The number of $j$ left of $i$ such that $P(j)$

# C-RASP: Example Program for Dyck-1

$$
\begin{array}{rcl}
Q_\ell & & \\
Q_r & & \\
C_\ell(i) & := & \#\,[j \le i]\; Q_\ell(j) \\
C_r(i) & := & \#\,[j \le i]\; Q_r(j) \\
V(i) & := & C_\ell(i) < C_r(i) \\
C_V(i) & := & \#\,[j \le i]\; V(j) \\
M(i) & := & C_V(i) = 0 \\
B(i) & := & C_\ell(i) = C_r \\
D(i) & := & M(i) \wedge B(i)
\end{array}
$$

# C-RASP Program Trace

## Program Trace

| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
|---|---|---|---|---|---|---|
| $Q_\ell$ | | | | | | |
| $Q_r$ | | | | | | |
| $C_\ell$ | | | | | | |
| $C_r$ | | | | | | |
| $V$ | | | | | | |
| $C_V$ | | | | | | |
| $M$ | | | | | | |
| $B$ | | | | | | |
| $D$ | | | | | | $\square$ |

## Initial Vectors - C-RASP Example: Dyck-$1$

The initial vectors are $Q_\ell$ and $Q_r$. These are all defined such that:

$Q_\ell(i)$ = True iff $\ell$ is $i$-th symbol

$Q_r(i)$ = True iff $r$ is $i$-th symbol

| | Program Trace | | | | | |
|---|---|---|---|---|---|---|
| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
| $Q_\ell$ | T | T | F | T | F | F |
| $Q_r$ | F | F | T | F | T | T |

$C_\ell$ counts the number of $\ell$ seen up until and including current position

$$C_\ell(i) = \#\,[j \le i]\ \ Q_\ell(i)\,.$$

| Program Trace | | | | | | |
|---|---|---|---|---|---|---|
| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
| $Q_\ell$ | T | T | F | T | F | F |
| $Q_r$ | F | F | T | F | T | T |
| $C_\ell$ | 1 | 2 | 2 | 3 | 3 | 3 |

$C_r$ counts the number of $r$ seen up until and including current position

$$C_r(i) = \# [j \leq i] \ Q_r(i).$$

| Program Trace | | | | | | |
|---|---|---|---|---|---|---|
| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
| $Q_\ell$ | T | T | F | T | F | F |
| $Q_r$ | F | F | T | F | T | T |
| $C_\ell$ | 1 | 2 | 2 | 3 | 3 | 3 |
| $C_r$ | 0 | 0 | 1 | 1 | 2 | 3 |

V indicates a matching violation - if there are ever more $r$ than $\ell$

$$V(i) = C_\ell(i) < C_r(i).$$

Program Trace

| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
|-------|--------|--------|-----|--------|-----|-----|
| $Q_\ell$ | T | T | F | T | F | F |
| $Q_r$ | F | F | T | F | T | T |
| $C_\ell$ | 1 | 2 | 2 | 3 | 3 | 3 |
| $C_r$ | 0 | 0 | 1 | 1 | 2 | 3 |
| $V$ | F | F | F | F | F | F |

# $C_V$ - C-RASP Example: Dyck$-1$

$C_V$ counts the number of violations seen up until the current position.

$$C_V(i) = \#\,[j \le i]\ \text{V(i)}.$$

| Program Trace | | | | | | |
|---|---|---|---|---|---|---|
| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
| $Q_\ell$ | T | T | F | T | F | F |
| $Q_r$ | F | F | T | F | T | T |
| $C_\ell$ | 1 | 2 | 2 | 3 | 3 | 3 |
| $C_r$ | 0 | 0 | 1 | 1 | 2 | 3 |
| $V$ | F | F | F | F | F | F |
| $C_V$ | 0 | 0 | 0 | 0 | 0 | 0 |

13

$M$ checks that the parentheses are always matched by verifying if there are zero violations

$$M(i) = C_V(i) = 0.$$

| | Program Trace | | | | | |
|---|---|---|---|---|---|---|
| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
| $Q_\ell$ | T | T | F | T | F | F |
| $Q_r$ | F | F | T | F | T | T |
| $C_\ell$ | 1 | 2 | 2 | 3 | 3 | 3 |
| $C_r$ | 0 | 0 | 1 | 1 | 2 | 3 |
| $V$ | F | F | F | F | F | F |
| $C_V$ | 0 | 0 | 0 | 0 | 0 | 0 |
| $M$ | T | T | T | T | T | T |

14

*B* checks that the parentheses are balanced by verifying if there are equally as many $\ell$ as $r$

$$B(i) = C_\ell(i) = C_r(i)$$

| | | Program Trace | | | | |
|---|---|---|---|---|---|---|
| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
| $Q_\ell$ | T | T | F | T | F | F |
| $Q_r$ | F | F | T | F | T | T |
| $C_\ell$ | 1 | 2 | 2 | 3 | 3 | 3 |
| $C_r$ | 0 | 0 | 1 | 1 | 2 | 3 |
| $V$ | F | F | F | F | F | F |
| $C_V$ | 0 | 0 | 0 | 0 | 0 | 0 |
| $M$ | T | T | T | T | T | T |
| $B$ | F | F | F | F | F | T |

D checks the string is matched and balanced.

$$D(i) = M(i) \wedge B(i)$$

Program Trace

| input | $\ell$ | $\ell$ | $r$ | $\ell$ | $r$ | $r$ |
|-------|--------|--------|-----|--------|-----|-----|
| $Q_\ell$ | T | T | F | T | F | F |
| $Q_r$ | F | F | T | F | T | T |
| $C_\ell$ | 1 | 2 | 2 | 3 | 3 | 3 |
| $C_r$ | 0 | 0 | 1 | 1 | 2 | 3 |
| $V$ | F | F | F | F | F | F |
| $C_V$ | 0 | 0 | 0 | 0 | 0 | 0 |
| $M$ | T | T | T | T | T | T |
| $B$ | F | F | F | F | F | T |
| $D$ | F | F | F | F | F | T |

# More Details

- Every **C-RASP** program compiles into a transformer that simulates it perfectly for inputs of arbitrary length
- Allows us to show the expressivity of transformers on many tasks: Dyck-1, $a^n b^n c^n$, and piecewise testable languages.

## What is in C-RASP and what isn't?

### $\in$ **C-RASP**

- DYCK-1
- Majority
- $a^n b^n c^n$
- Piecewise testable
  $\Sigma^* a_1 \Sigma^* a_2 \Sigma^* \dots \Sigma^* a_n \Sigma^*$

### $\notin$ **C-RASP**
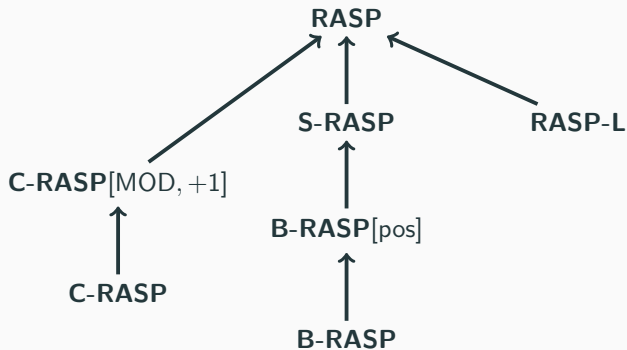
- $\Sigma^* a c^* a \Sigma^*$
- $\{a^m \mid m = n^2, n \in \mathbb{N}\}$
- $\{w\$w \mid \text{for } w \in \Sigma\}$
- $(aa)^*$ and PARITY
- **NC$^1$**-complete languages

some of these are just conjectures . . .

## The RASP Family Tree



Disclaimer: More arrows may exist

- Every **C-RASP** program compiles into a transformer that simulates it perfectly for inputs of arbitrary length

- Understand what transformer encoders can and cannot do, more intuitively

- Connections with formal language theory and logic

Thanks!

Dan Friedman, Alexander Wettig, and Danqi Chen. Learning transformer programs. *Advances in Neural Information Processing Systems*, 36, 2024.

David Lindner, János Kramár, Sebastian Farquhar, Matthew Rahtz, Tom McGrath, and Vladimir Mikulik. Tracr: Compiled transformers as a laboratory for interpretability. *Advances in Neural Information Processing Systems*, 36, 2024.

Lena Strobl. Average-hard attention transformers are constant-depth uniform threshold circuits, 2023. URL https://arxiv.org/abs/2308.03212. arXiv:2308.03212.

Lena Strobl, William Merrill, Gail Weiss, David Chiang, and Dana Angluin. Transformers as recognizers of formal languages: A survey on expressivity. *arXiv preprint arXiv:2311.00208*, 2023.

Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International Conference on Machine Learning*, pages 11080–11090. PMLR, 2021.

Andy Yang and David Chiang. Counting like transformers: Compiling temporal counting logic into softmax transformers. In *Proceedings of the Conference on Language Modeling*, 2024. URL https://arxiv.org/abs/2404.04393. To appear.

# Extra C-RASP Examples

# More?

$a^* b^*$ over $\Sigma = \{a, b\}$

$$C_a(i) := \#\,[j \leq i]\;Q_a(j)$$      # positions with $a$'s

$$C_b(i) := \#\,[j \leq i]\;Q_b(j)$$      # positions with $b$'s

$$V(i) := Q_a(i) \land C_b(i) \geq 1$$      *Violation*: an $a$ has $b$'s preceding it

$$C_V(i) := \#\,[j \leq i]\;V(j)$$      # *Violations*

$$Y(i) := C_V(i) = 0$$      Zero *Violations*

22

$a^* b^* a^*$ over $\Sigma = \{a, b\}$

$$C_a(i) := \#\,[j \leq i]\ \ Q_a(j) \qquad\qquad \text{\# positions with } a\text{'s}$$

$$C_b(i) := \#\,[j \leq i]\ \ Q_b(j) \qquad\qquad \text{\# positions with } b\text{'s}$$

$$BA(i) := Q_a(i) \wedge C_b(i) \geq 1 \qquad\qquad \text{A subsequence } ba \text{ ends at } i$$

$$C_{ba}(i) := \#\,[j \leq i]\ \ BA(j) \qquad\qquad \text{\# ends of subsequence } ba$$

$$BAB(i) := Q_b(i) \wedge C_{ba} \geq 1 \qquad\qquad \text{the subsequence } bab \text{ ends at } i$$

$$C_{bab}(i) := \#\,[j \leq i]\ \ BAB(j) \qquad\qquad \text{\# ends of subsequence } bab$$

$$Y(i) := C_{bab}(i) = 0 \qquad\qquad \text{There are no subsequences } bab$$

# MORE?!

$a^n b^n c^n$ over $\Sigma = \{a, b, c\}$

$$C_a(i) := \#\,[j \leq i]\ Q_a(j) \qquad\qquad \text{\# positions with } a\text{'s}$$
$$C_b(i) := \#\,[j \leq i]\ Q_b(j) \qquad\qquad \text{\# positions with } b\text{'s}$$
$$C_c(i) := \#\,[j \leq i]\ Q_c(j) \qquad\qquad \text{\# positions with } c\text{'s}$$
$$A(i) := C_b(i) + C_c(i) = 0 \qquad\qquad \text{No preceding } b\text{'s or } c\text{'s}$$
$$B(i) := C_c(i) = 0 \qquad\qquad \text{No preceding } c\text{'s}$$
$$C_A(i) := \#\,[j \leq i]\ Q_a(j) \wedge A(j) \qquad\qquad \text{\# } a\text{'s with no preceding } b\text{'s or } c\text{'s}$$
$$C_B(i) := \#\,[j \leq i]\ Q_b(j) \wedge B(j) \qquad\qquad \text{\# } b\text{'s with no preceding } c\text{'s}$$
$$G_a(i) := C_A(i) = C_a(i) \qquad\qquad \text{no } a\text{'s have preceding } b\text{'s or } c\text{'s}$$
$$G_b(i) := C_B(i) = C_b(i) \qquad\qquad \text{no } b\text{'s have preceding } c\text{'s}$$
$$G_{abc}(i) := C_a(i) = C_b(i) = C_c(i) \qquad\qquad \text{same number of } a\text{'s, } b\text{'s, } c\text{'s}$$
$$Y(i) := G_a(i) \wedge G_b(i) \wedge G_{abc}(i) \qquad\qquad \text{Correct order \& number of symbols}$$

# Ok one more

*hello over* $\Sigma = \{e, h, l, o\}$

$$C_e(i) := \#\,[j \leq i]\;\; Q_e(j) \qquad\qquad \text{\# positions with } e\text{'s}$$
$$C_h(i) := \#\,[j \leq i]\;\; Q_h(j) \qquad\qquad \text{\# positions with } h\text{'s}$$
$$C_l(i) := \#\,[j \leq i]\;\; Q_l(j) \qquad\qquad \text{\# positions with } l\text{'s}$$
$$C_o(i) := \#\,[j \leq i]\;\; Q_o(j) \qquad\qquad \text{\# positions with } o\text{'s}$$
$$C_\Sigma(i) := \#\,[j \leq i]\;\; 1 \qquad\qquad \text{\# symbols in string}$$
$$HE(i) := Q_e(i) \land C_h(i) = 1 \qquad\qquad \text{A subsequence } he \text{ ends at } i$$
$$C_{he}(i) := \#\,[j \leq i]\;\; HE(j) \qquad\qquad \text{\# ends of subsequence } he$$
$$HEL(i) := Q_l(i) \land C_{he}(i) = 1 \qquad\qquad \text{A subsequence } hel \text{ ends at } i$$
$$C_{hel}(i) := \#\,[j \leq i]\;\; HEL(j) \qquad\qquad \text{\# ends of subsequence } hel$$
$$HELLO(i) := Q_o(i) \land C_{hel}(i) = 2 \qquad\qquad \text{A subsequence } hello \text{ ends at } i$$
$$Y(i) := HELLO(i) \land C_\Sigma(i) = 5 \qquad\qquad \text{Length 5 and contains } hello$$