

A Formal Framework For Length Generalization in Transformers

Xinting Huang^{1*} Andy Yang^{2*} Satwik Bhattamishra³ Yash Sarrof¹
Andreas Krebs⁴ Hattie Zhou⁵ Preetum Nakkiran⁶ Michael Hahn¹

¹Saarland University

²University of Notre Dame

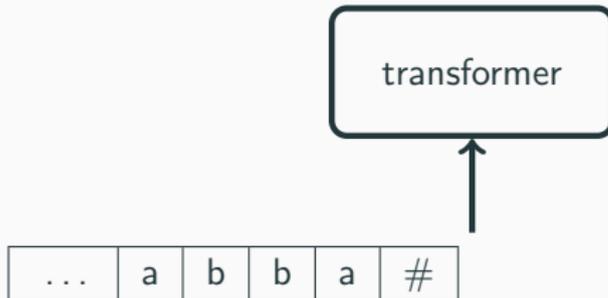
³University of Oxford

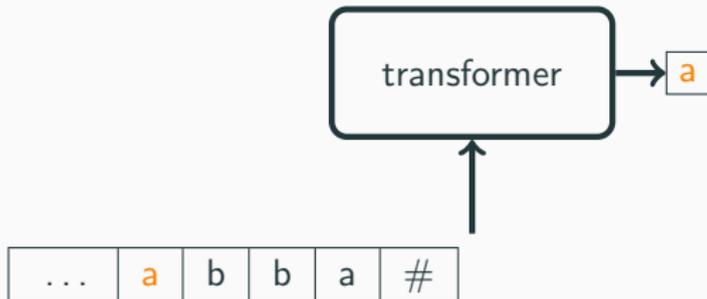
⁴University of Tübingen

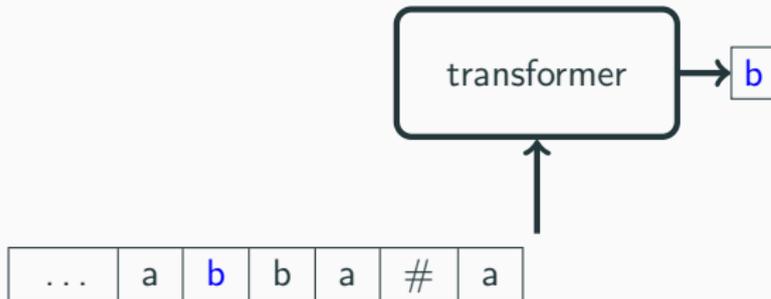
⁵Mila, Université de Montréal

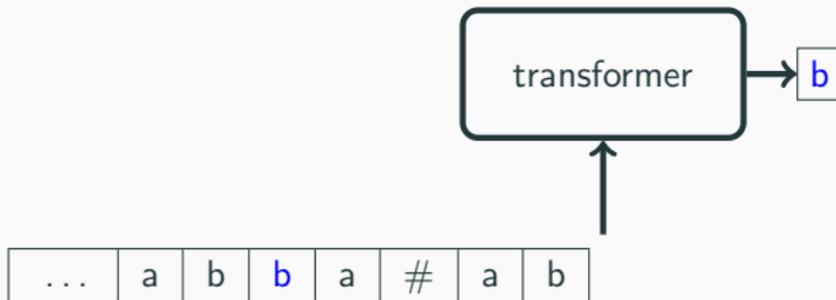
⁶Apple

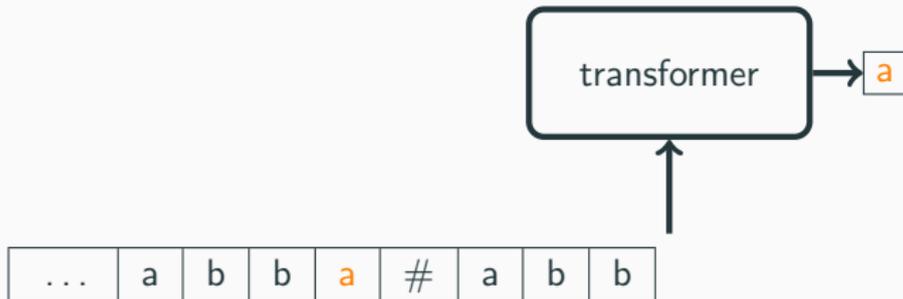


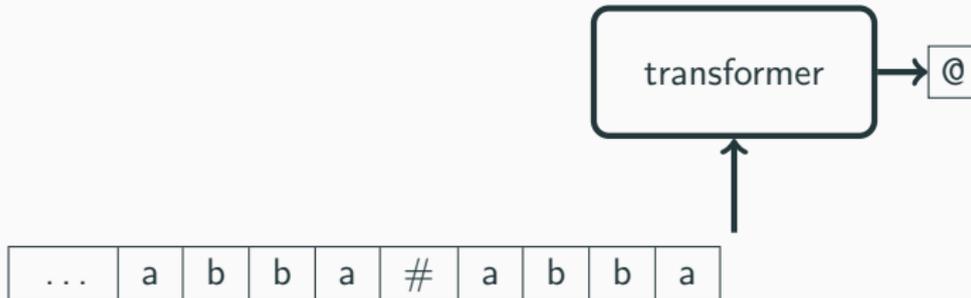


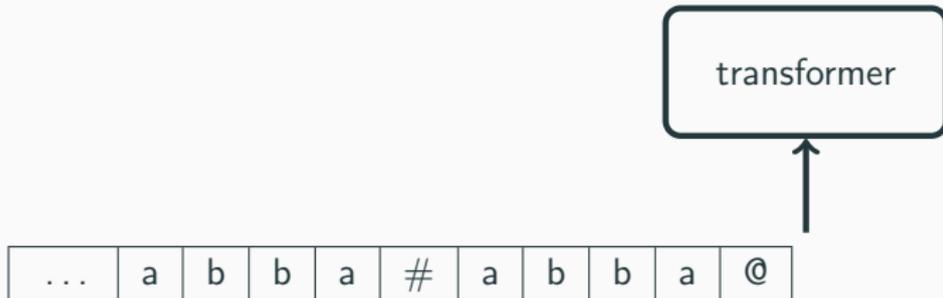












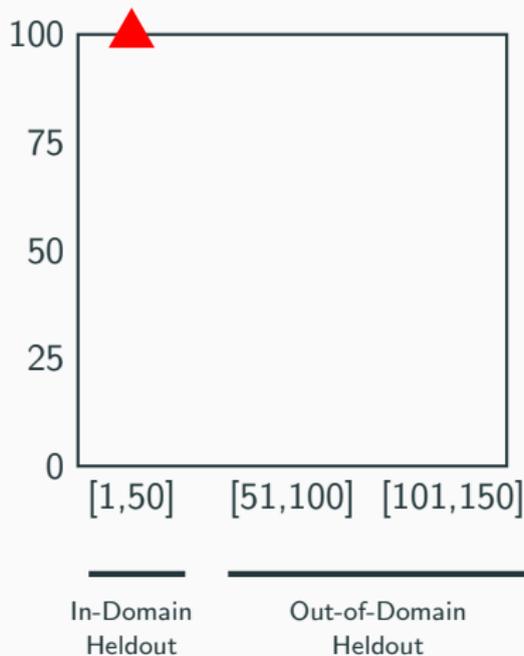
Length Generalizing on Copy

If we train on length up to $N \dots$



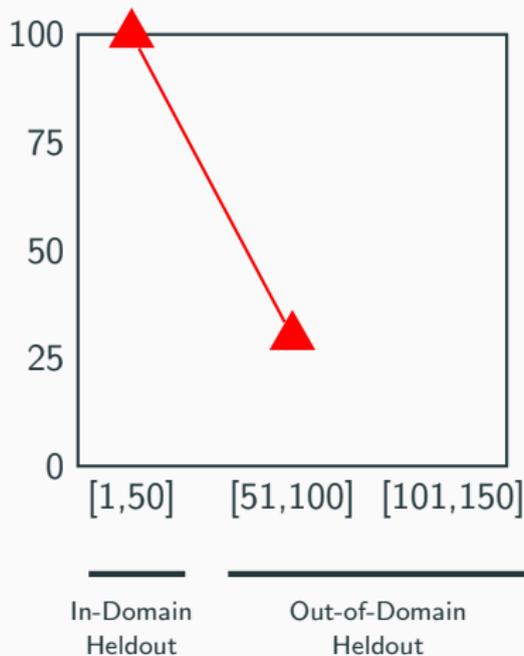
Will we succeed on lengths up to $2N \dots$?

Length Generalizing on Copy



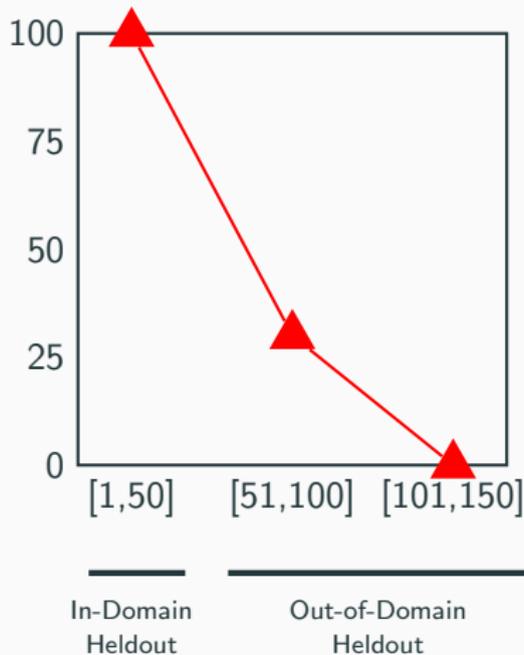
Zhou et al 2024 ICLR; Jelassi et al 2024 ICML; Kazemnejad et al 2023 NeurIPS.

Length Generalizing on Copy



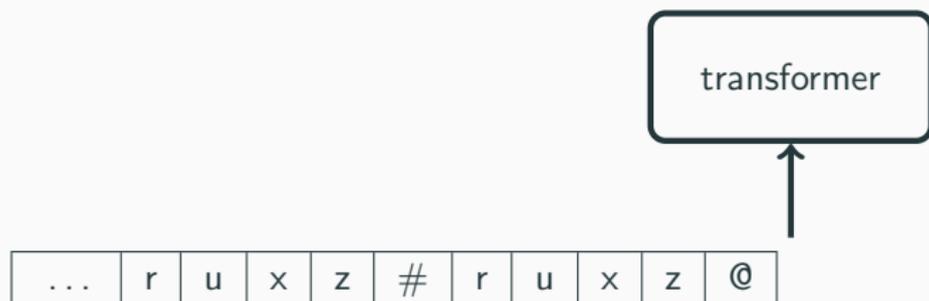
Zhou et al 2024 ICLR; Jelassi et al 2024 ICML; Kazemnejad et al 2023 NeurIPS.

Length Generalizing on Copy

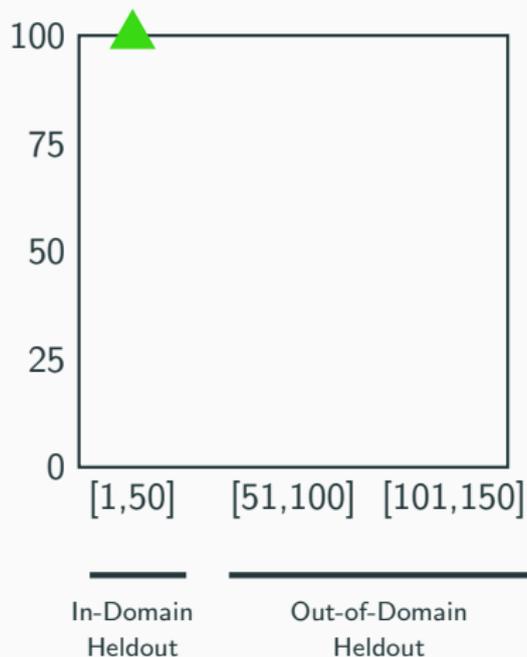


Zhou et al 2024 ICLR; Jelassi et al 2024 ICML; Kazemnejad et al 2023 NeurIPS.

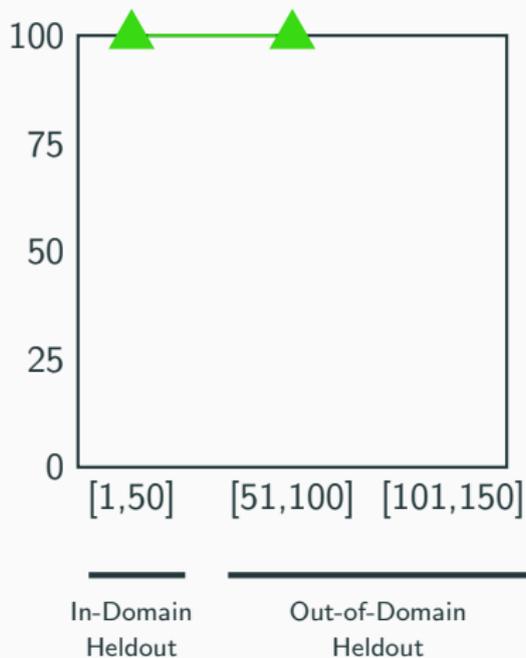
Unique Copy



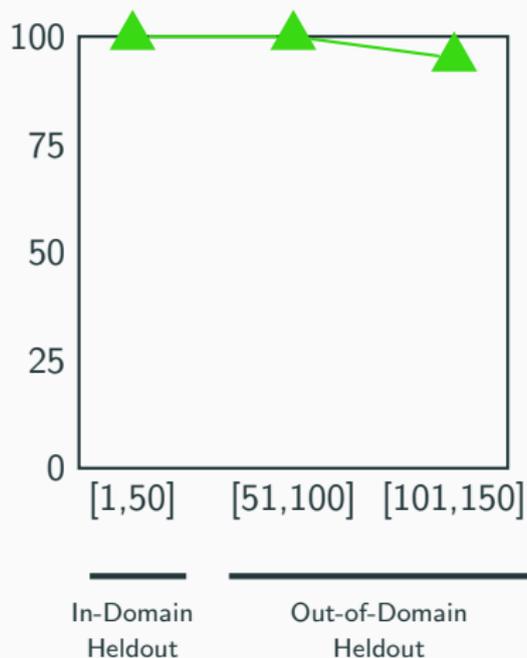
Length Generalizing on Unique Copy



Length Generalizing on Unique Copy



Length Generalizing on Unique Copy



On which tasks do transformers
length-generalize?

What Algorithms can Transformers Learn? A Study in Length Generalization

Hattie Zhou Arwen Bradley Etai Littwin Noam Razin
Omid Saremi, Josh Susskind Samy Bengio Preetum Nakkiran

RASP-L Conjecture (paraphrased): Transformers length-generalize on problems with simple programs in the RASP-L language.

RASP-L Conjecture (paraphrased): Transformers length-generalize on problems with simple programs in the RASP-L language.

Challenge: RASP-L hasn't been formalized. Expressiveness not understood

RASP-L Conjecture (paraphrased): Transformers length-generalize on problems with simple programs in the RASP-L language.

Challenge: RASP-L hasn't been formalized. Expressiveness not understood

Our Contribution: Formalize it based on C-RASP.

MAJORITY

$C_1(i) := \# [j \leq i] Q_1(j)$

count # of 1's (1)

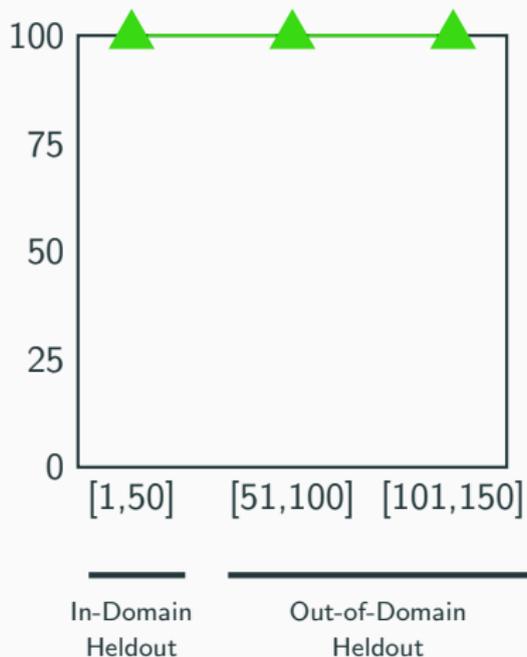
$C_0(i) := \# [j \leq i] Q_0(j)$

count # of 0's (2)

$M(i) := C_1(i) \geq C_0(i)$

compare them (3)

Length Generalizing on Majority



C-RASP Example: $(aa)^*$

$(aa)^*$

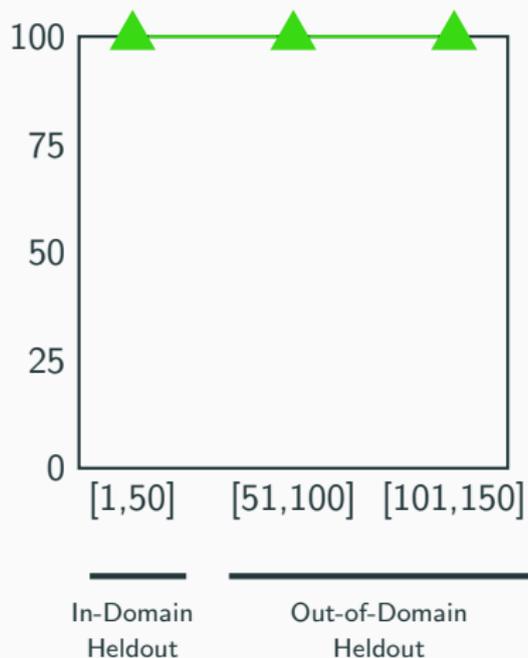
$\phi(i) := i \equiv 0 \pmod{2}$ check position modulo 2 (1)

$C_{\neg a}(i) := \# [j \leq i] \neg Q_a(j)$ (2)

$A(i) := C_{\neg a}(i) = 0$ check that no b 's occur (3)

$D(i) := \phi(i) \wedge A(i)$ (4)

Length Generalizing on $(aa)^*$



C-RASP Example: Unique Copy

Unique Copy (Induction Head)

\vdots (1)

$CP_a(i) := \# [j \leq i, j = i - 1] Q_a(j)$ (2)

$PRED_a(i) := CP_a(i) \geq 1$ check preceding token (3)

\vdots (4)

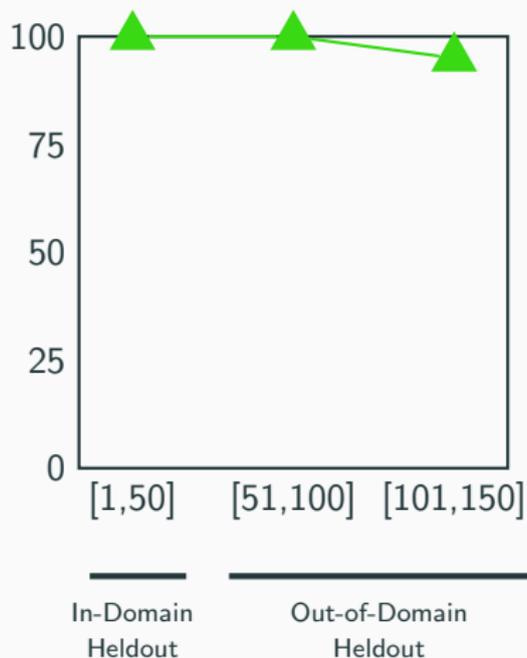
$CBIGRAM_{ab} := \# [j \leq i] Q_b(j) \wedge PRED_a(j)$ collect bigrams (5)

$EXISTS_{ab} := CBIGRAM_{ab}(i) \geq 1$ (6)

\vdots (7)

$NEXT_a(i) := \bigvee_{\sigma \in \Sigma} [Q_\sigma(i) \wedge EXISTS_{\sigma a}(i)]$ check if bigram matches (8)

Length Generalizing on Unique Copy



Boolean-Valued Operations

Initial $P(i) := Q_\sigma(i)$
 for $\sigma \in \Sigma$

Boolean $P(i) := \neg P_1(i)$
 $P(i) := P_1(i) \wedge P_2(i)$

Constant $P(i) := \top$

Positional $P(i) := \phi(i)$
 for $\phi \in \Phi$

Comparison $P(i) := C_1(i) \leq C_2(i)$

Count-Valued Operations

Counting $C(i) := \# [j \leq i, \psi(i, j)] P(j)$
 for $\psi \in \Psi \cup \{\top\}$

Conditional $C(i) := P(i) ? C_1(i) : C_2(i)$

Addition $C(i) := C_1(i) + C_2(i)$

Subtraction $C(i) := C_1(i) - C_2(i)$

Constant $C(i) := 1$

Copy With Repetition

⋮	(1)
⋮	(2)
???	(3)

Provably no C-RASP program!

Rigorous proof via communication complexity.

Theorem (informal):

If f has a C-RASP program, then transformers trained on an idealized training procedure will length-generalize from N to $2N$, when N is large.

Definition: Inference Procedure

Let f be the target function. For $n = 1, 2, 3, \dots$ choose a transformer T_n

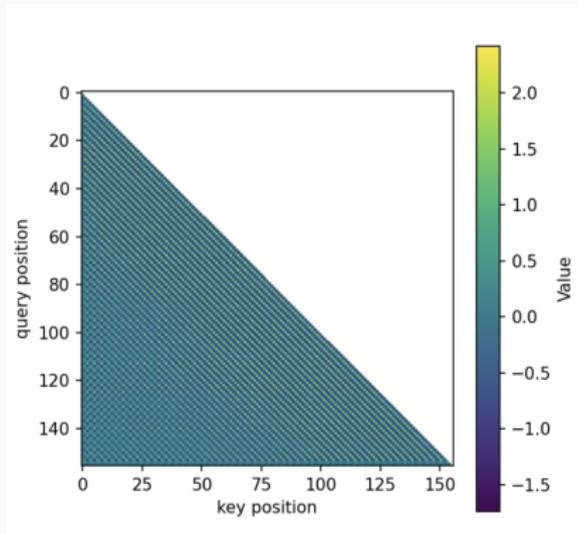
- that has context window n
- matches f on all input strings with length $\leq \frac{n}{2}$
- while minimizing a regularizer

This procedure converges on a length-generalizing transformer. There is some N such that, for $n > N$, T_n matches f on all strings up to length n .

Translation Invariance

Our hypothesis class assumes translation invariance in transformers, following the APE set up.

While not enforced by SGD, we argue translation invariance is helpful for transformers to length generalize

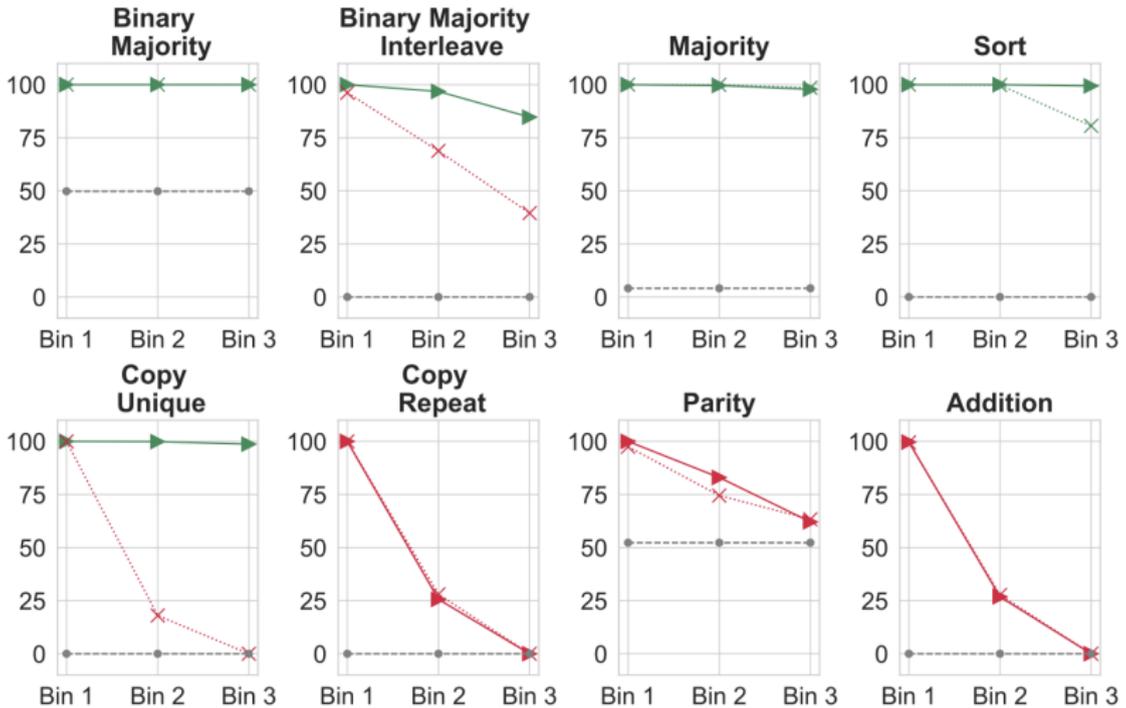


Definition (Regularizer)

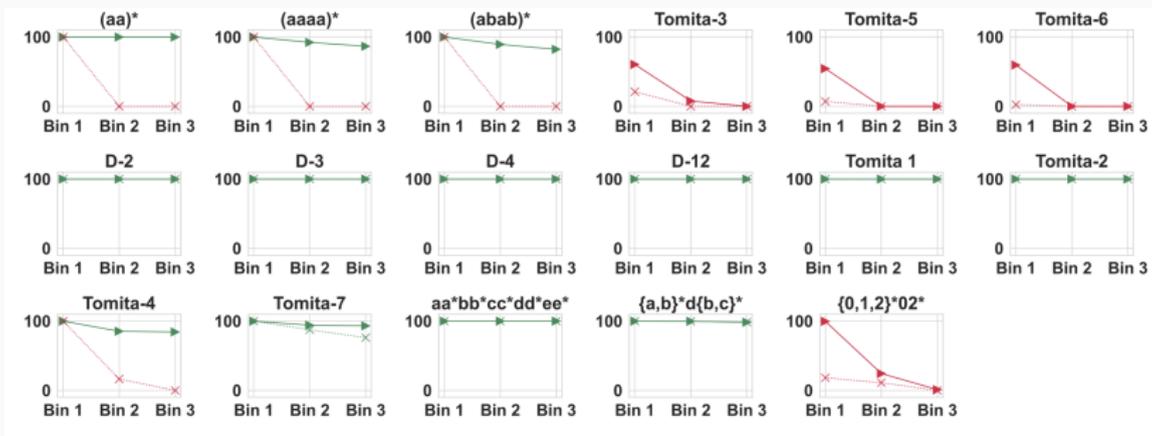
Define $\mathcal{R}(T)$ as the sum of (1) $L + H$; (2) the precision p used in the hypothesis class; the precision p used for rounding attention logits and the output of $\exp(\cdot)$; (3) $\max_{l,h} \text{rank}(\mathbf{V}_{l,h})$; (4) $\max_{l,h} \|\mathbf{K}_{l,h}^T \mathbf{Q}_{l,h}\|$; $\max_{l,h} \|\mathbf{V}_{l,h}\|$; $\max_l \|\mathbf{A}_l\|_F$, $\|\mathbf{B}_l\|_F$; $\|\mathbf{U}\|$; (5) $\max_i \|\mathbf{p}_i\|_2$, $\max_\sigma \|\mathbf{E}_\sigma\|_2$, $\max_l \|\mathbf{b}_l\|_2$; (6) the term

$$\sum_{l=1}^L \sum_{h=1}^H \sum_{j=1}^n |\mathbf{p}_1^T \mathbf{K}_{l,h}^T \mathbf{Q}_{l,h} \mathbf{p}_j|^2 \quad (4)$$

Algorithmic Tasks



Formal Language Tasks



Bhattacharya et al. 2020

On which tasks do transformers length-generalize?

Theorem (informal): If f has a C-RASP program, then transformers trained on an idealized training procedure will length-generalize from N to $2N$, when N is large.

MAJORITY

$$C_1(i) := \#\{j \leq i \mid Q_1(j)\} \quad (1)$$

$$C_0(i) := \#\{j \leq i \mid Q_0(j)\} \quad (2)$$

$$M(i) := C_1(i) \geq C_0(i) \quad (3)$$

