



Evaluating Transformer's Ability to Learn Mildly Context-Sensitive Languages

Shunju Wang, Shane Steinert-Threlkeld

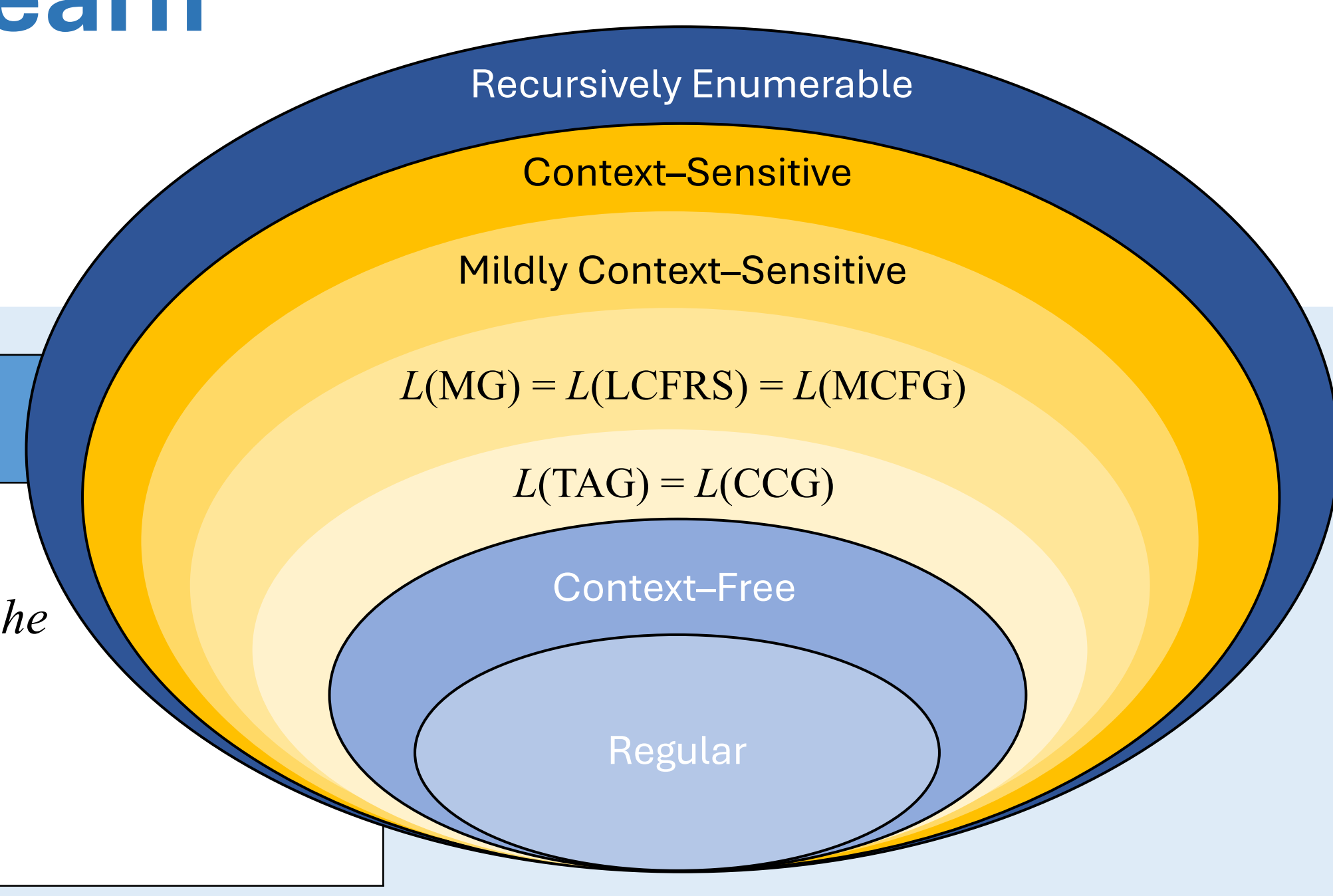


Natural Language is Supra-Context-Free

Swiss German subordinate clauses cross-serial dependency (Shieber, 1985)

Jan säit das mer d'chind em Hans es huus haend wele laa hülfe aastriiche
 Jan says that we the children-ACC Hans-DAT the house-ACC have wanted let help paint

'Jan says that we have wanted to let the children help Hans paint the house.'



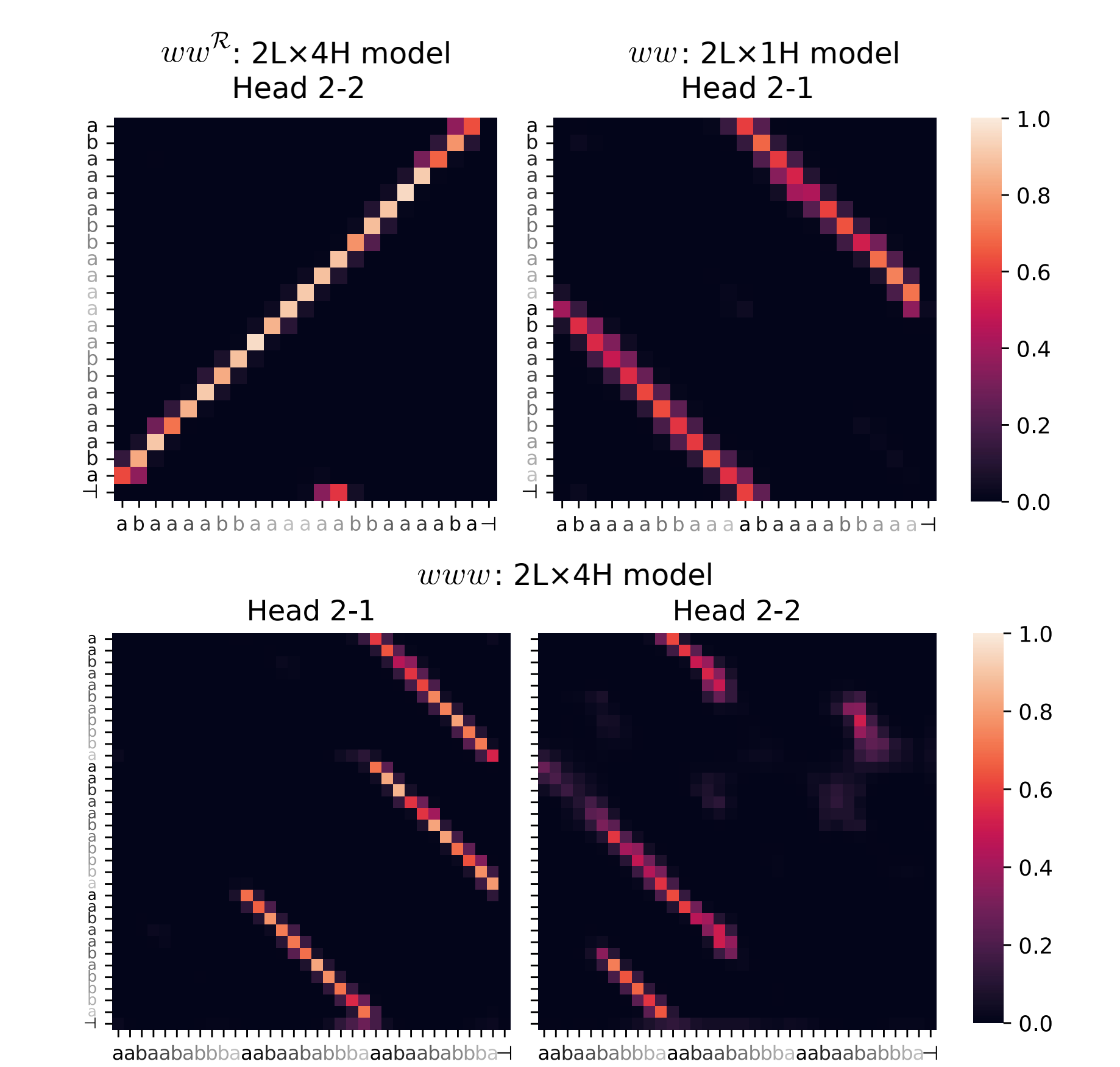
Mild Context-Sensitivity

- TAG, CCG, etc. extend CFG with just enough power to describe cross-serial dependency as in Swiss German.
- MG, MCFG, etc. have power beyond TAG as motivated by more complex phenomena.
- Formalization of *mildly context-sensitive languages* (Kallmayer 2010):
 - Describe cross-serial dependencies
 - Can be parsed in polynomial time
 - String length grows linearly
 - Contain all CFLs
- MCSLs as benchmarks for linguistic adequacy:
 - Represent a hypothesized upper bound of the complexity of natural language
 - Abstractions of complex phenomena such as reduplication, free word order, etc.

Copying

BINARY CLASSIFICATION
 POS: $\{ww \mid w \in \{a, b\}^*\}$
 NEG: random strings from $\{a, b\}^*$

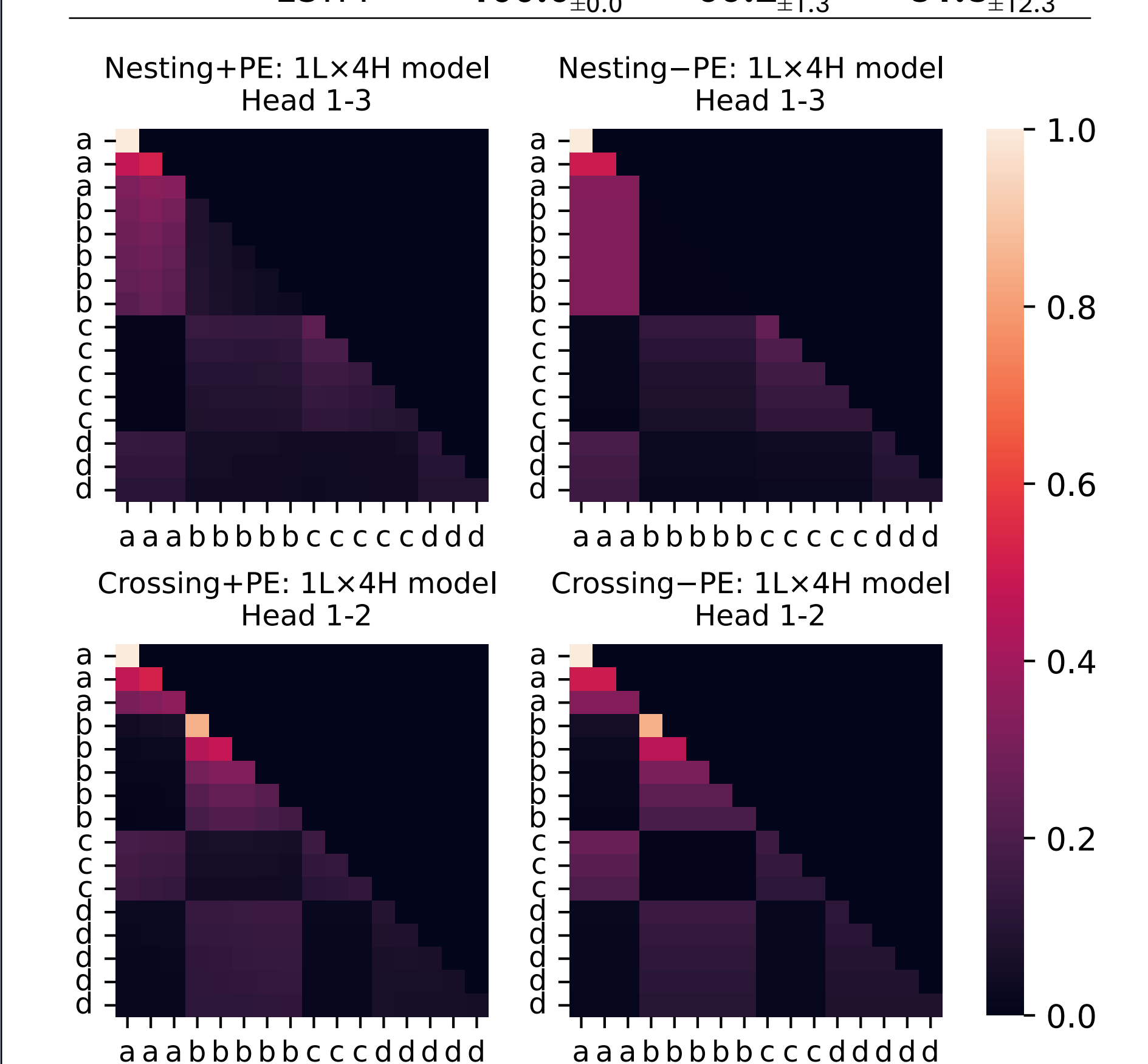
Accuracy (%)	IN-DISTR.	OOD	OOD	
	$ w \in [1, 11]$	$ w = 12$	$ w = 13$	
ww^R	Transf.	99.5 \pm 0.3	50.4 \pm 0.3	50.2 \pm 0.1
	LSTM	97.8 \pm 0.5	96.0 \pm 0.7	96.0 \pm 0.8
ww	Transf.	99.5 \pm 0.1	51.3 \pm 0.3	50.2 \pm 0.0
	LSTM	97.2 \pm 0.4	95.7 \pm 1.1	90.4 \pm 1.4
www	Transf.	99.5 \pm 0.2	51.0 \pm 0.5	50.5 \pm 0.2
	LSTM	99.4 \pm 0.1	98.6 \pm 0.5	87.5 \pm 6.2



Crossing

NEXT CHARACTER PREDICTION
 $a^n b^m c^n d^m \rightarrow (a/b)^n (b/c)^m c^{n-1} d^m$ [EOS]

Accuracy (%)	IN-DISTR.	OOD	OOD	
	$n, m \in [1, 50]$	$n \text{ or } m \in [51, 100]$	$n \text{ or } m \in [101, 150]$	
$a^n b^m c^m d^n$	Tr. +PE	99.8 \pm 0.2	6.5 \pm 1.3	0.0 \pm 0.0
	Tr. -PE	100.0 \pm 0.0	98.0 \pm 0.3	23.0 \pm 3.1
$a^n b^m c^n d^m$	LSTM	100.0 \pm 0.0	100.0 \pm 0.0	100.0 \pm 0.0
	Tr. +PE	100.0 \pm 0.0	7.2 \pm 1.6	0.0 \pm 0.0
$a^n b^m c^n d^m$	Tr. -PE	100.0 \pm 0.0	92.3 \pm 1.2	27.0 \pm 14.0
	LSTM	100.0 \pm 0.0	99.2 \pm 1.3	81.3 \pm 12.3



We test how well Transformers learn complex MCSLs. They generalize well

to unseen in-distribution data, but their extrapolation is worse than LSTMs. The learned self-attention resembles dependency relations and the representations encoded count information.

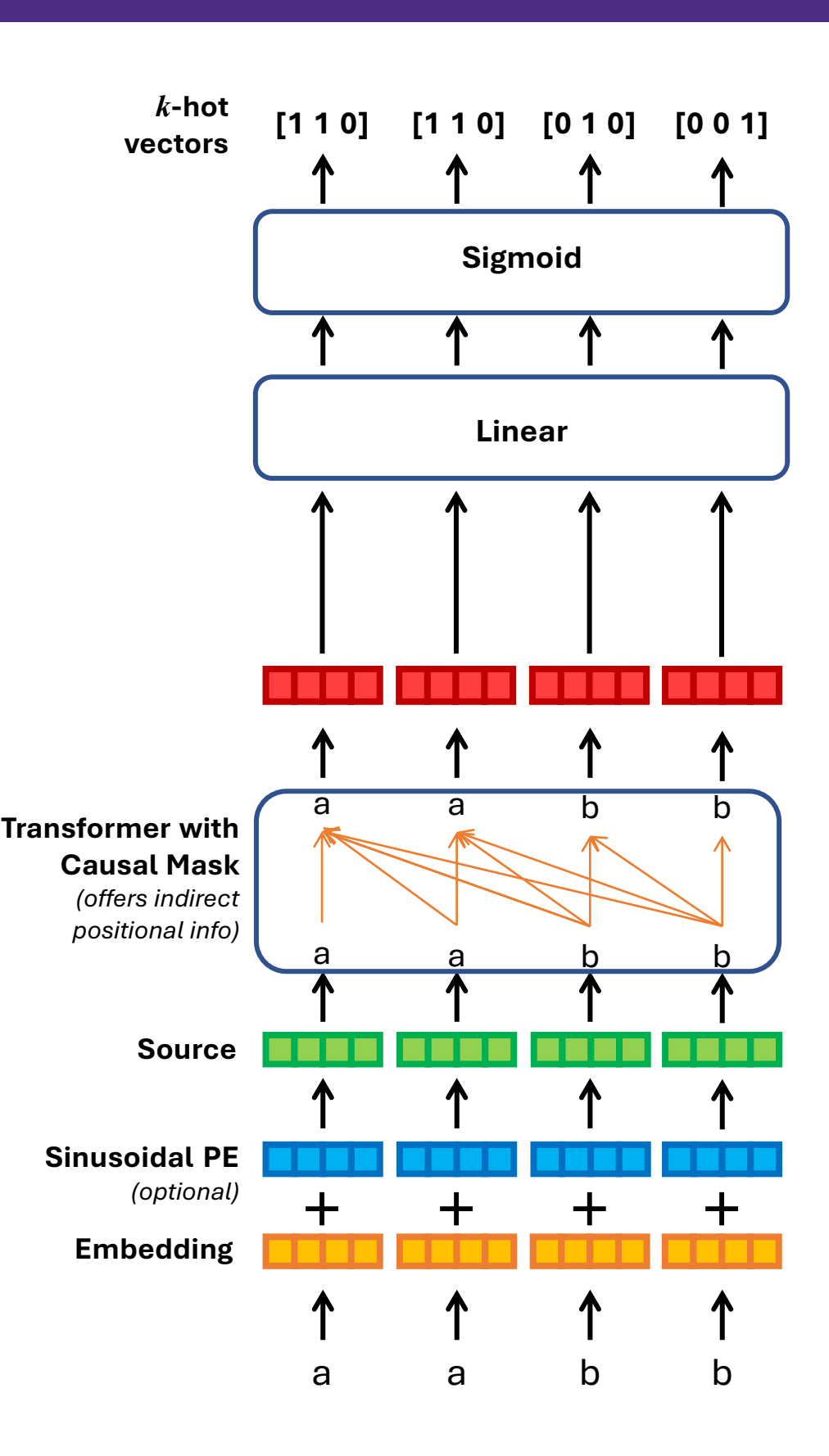
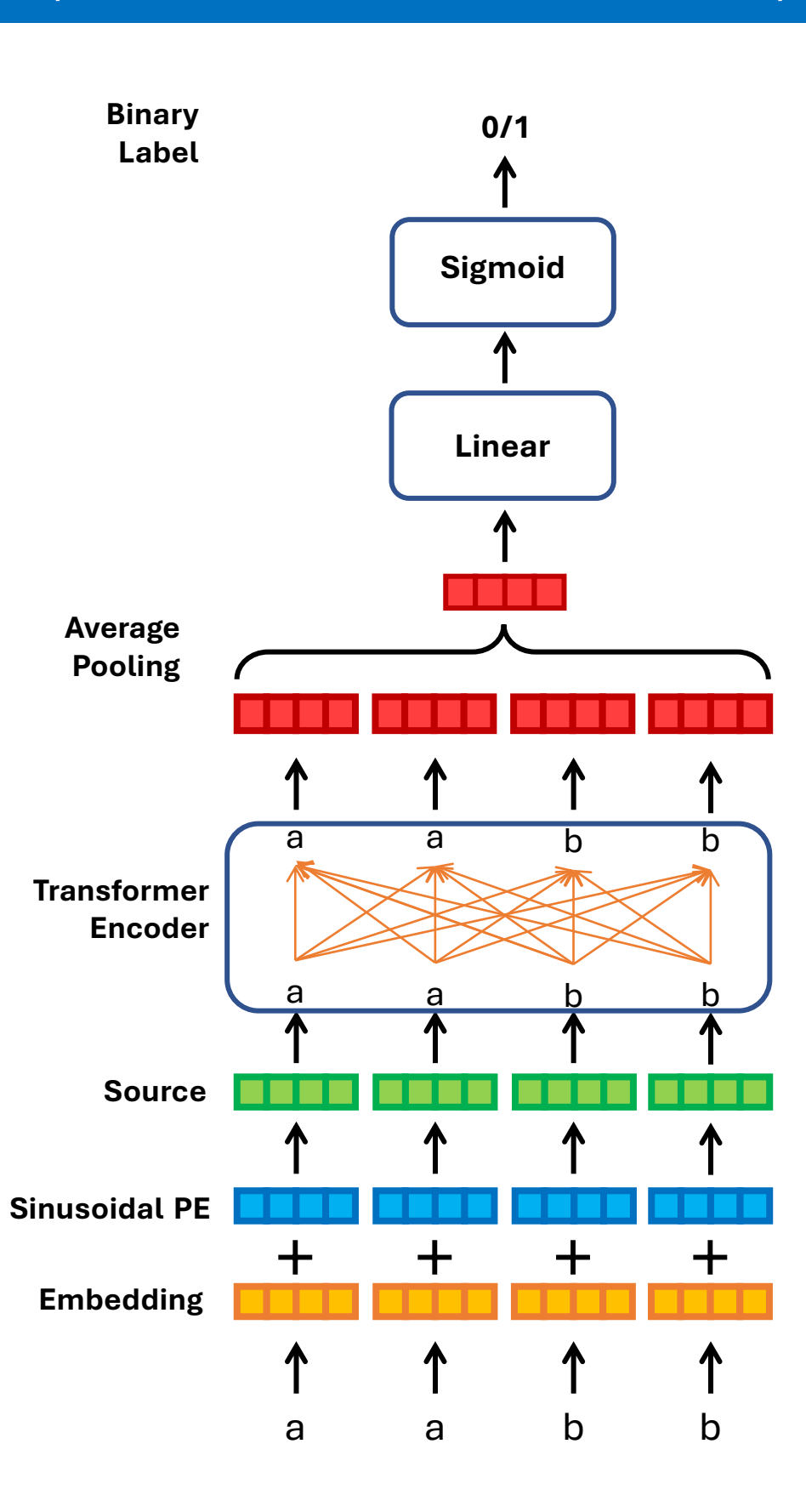
Languages

CFL	Mildly Context-Sensitive	
	$L(\text{TAG})$	$L(\text{MG}) = L(\text{MCFG})$
LESS COMPLEX	CANONICAL	MORE COMPLEX
		SCRAMBLE
ww^R	ww	www
$a^n b^m c^m d^n$	$a^n b^m c^n d^m$	$O_2 = \{w \in \{a, b, c, d\}^* \mid w _a = w _c \wedge w _b = w _d\}$
$a^n b^n$	$a^n b^n c^n$	$a^n b^n c^n d^n e^n$
	$a^n b^n c^n d^n$	$\text{MIX} = \{w \in \{a, b, c\}^* \mid w _a = w _b = w _c\}$

Tasks

BINARY CLASSIFICATION (BIDIRECTIONAL ATTENTION)

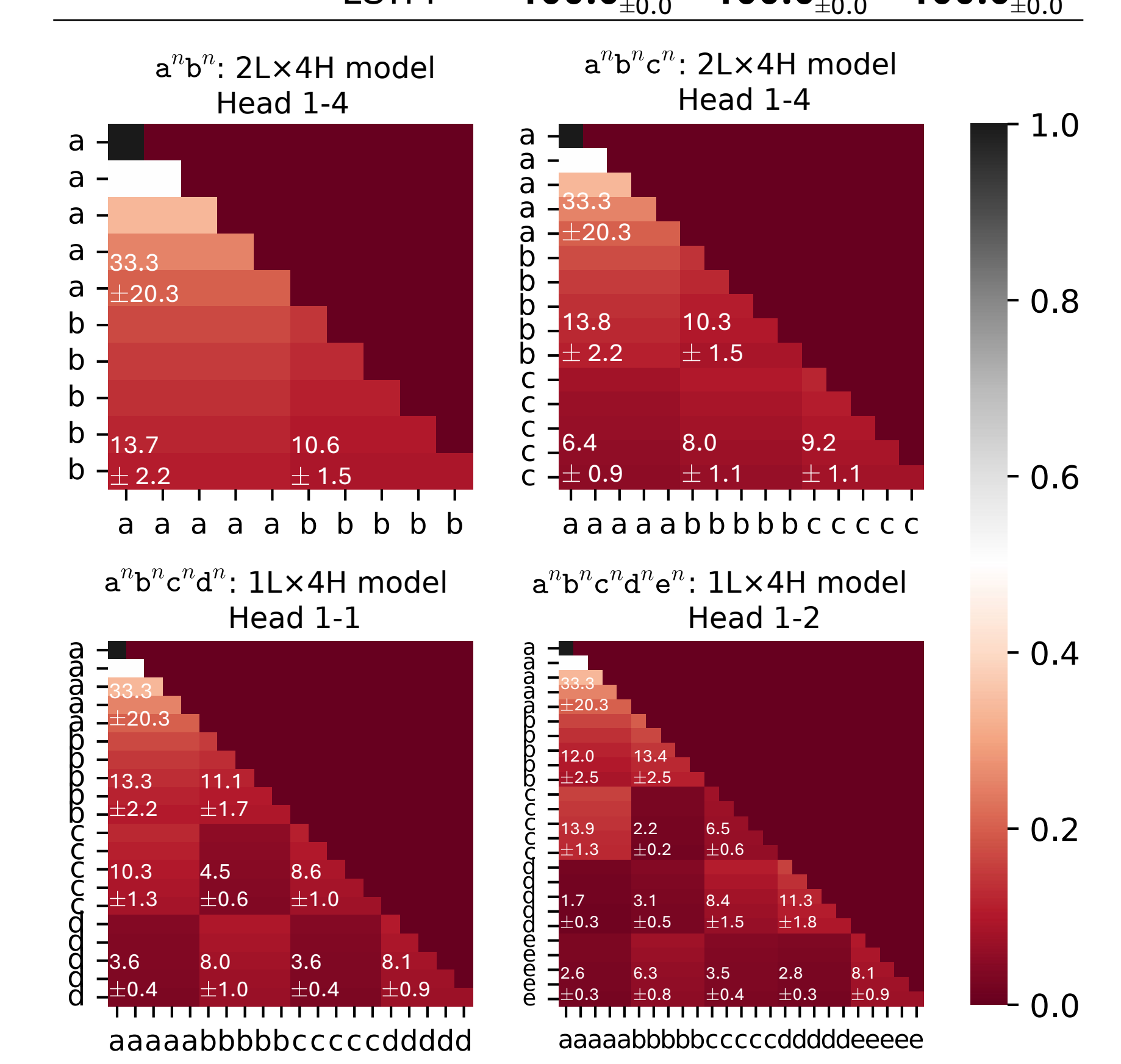
NEXT CHARACTER PREDICTION (UNIDIRECTIONAL ATTENTION)



Multiple Agreements

NEXT CHARACTER PREDICTION
 $a^n b^n c^n d^n e^n \rightarrow (a/b)^n b^{n-1} c^n d^n e^n$ [EOS]

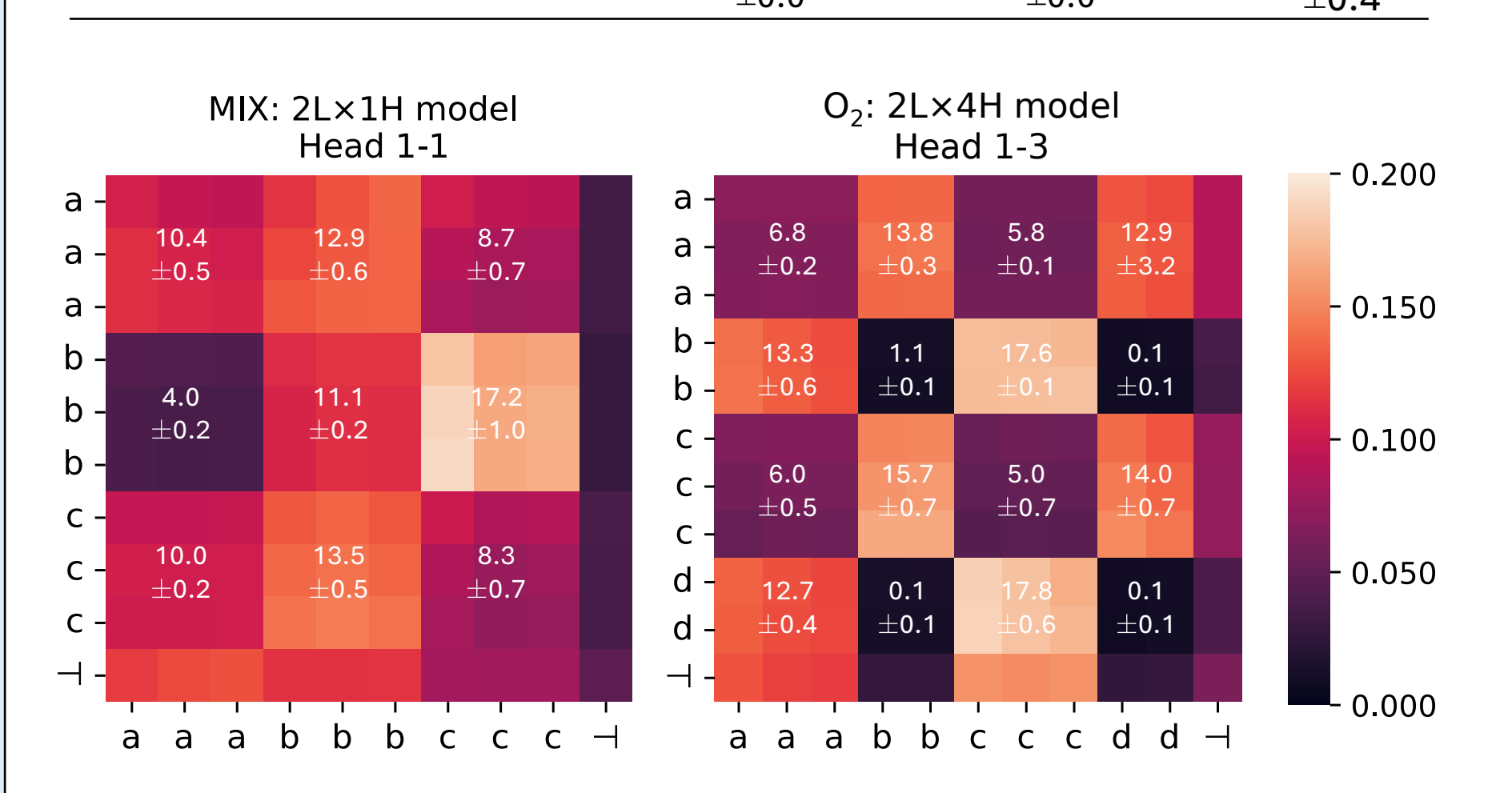
Accuracy (%)	IN-DISTR.	OOD	OOD	
	$n \in [1, 50]$	$n \in [51, 100]$	$n \in [101, 150]$	
$a^n b^n$	Tr. -PE	100.0 \pm 0.0	100.0 \pm 0.0	91.3 \pm 8.4
	LSTM	100.0 \pm 0.0	100.0 \pm 0.0	100.0 \pm 0.0
$a^n b^n c^n$	Tr. -PE	100.0 \pm 0.0	100.0 \pm 0.0	36.0 \pm 14.2
	LSTM	100.0 \pm 0.0	100.0 \pm 0.0	100.0 \pm 0.0
$a^n b^n c^n d^n$	Tr. -PE	100.0 \pm 0.0	100.0 \pm 0.0	24.0 \pm 10.2
	LSTM	100.0 \pm 0.0	100.0 \pm 0.0	48.7 \pm 13.6
$a^n b^n c^n d^n e^n$	Tr. -PE	100.0 \pm 0.0	85.3 \pm 15.4	3.3 \pm 4.7
	LSTM	100.0 \pm 0.0	100.0 \pm 0.0	100.0 \pm 0.0



Scrambling

BINARY CLASSIFICATION
 POS: all permutations of $a^n b^n c^n / a^n b^m c^n d^m$
 NEG: remaining strings from $\{a, b, c\}^* / \{a, b, c, d\}^*$

Macro F-1(%)	IN-DISTR.	OOD	OOD	
	$ w _\sigma \in [1, 4]$	$ w _\sigma = 5$	$ w _\sigma = 6$	
MIX	Transformer	100.0 \pm 0.0	65.6 \pm 2.9	45.7 \pm 6.3
	LSTM	100.0 \pm 0.0	70.3 \pm 10.5	49.0 \pm 15.5
O_2	Transformer	100.0 \pm 0.0	60.5 \pm 8.5	45.1 \pm 10.1
	LSTM	100.0 \pm 0.0	100.0 \pm 0.0	98.6 \pm 0.4



Using an MLP regressor prober, we can extract the ongoing tallies for the 3 symbols in MIX strings. The predictions and targets have an MSE of 0.21 and a Pearson correlation of 0.929. This contrasts with a control task target (shuffled original target) which has an MSE of 1.33.

Counting Target			
	#a	#b	#c
a	[1	0	0]
b	[1	1	0]
c	[1	1	1]
a	[2	1	1]
b	[2	2	1]
c	[2	2	2]