



#### **About Me**

- Postdoc at ETH Zürich with Ryan Cotterell
- Formerly: PhD student at University of Notre Dame with David Chiang
- Interests
  - Natural language processing
  - Formal language theory
  - Neural networks





#### **Basic Research Questions**

# What can LLMs do? What can't LLMs do?

Published as a conference paper at ICLR 2025

# TRAINING NEURAL NETWORKS AS RECOGNIZERS OF FORMAL LANGUAGES

Alexandra Butoi<sup>1</sup> Ghazal Khalighinejad<sup>2</sup> Anej Svete<sup>1</sup>

Josef Valvoda<sup>3</sup> Ryan Cotterell<sup>1</sup> Brian DuSell<sup>1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>Duke University <sup>3</sup>University of Copenhagen {alexandra.butoi, anej.svete, ryan.cotterell, brian.dusell}@inf.ethz.ch ghazal.khalighinejad@duke.edu jval@di.ku.dk

#### **ABSTRACT**

Characterizing the computational power of neural network architectures in terms of formal language theory remains a crucial line of research, as it describes lower and upper bounds on the reasoning capabilities of modern AI. However, when empirically testing these bounds, existing work often leaves a discrepancy between experiments and the formal claims they are meant to support. The problem is that

# **My Collaborators**



Alexandra Butoi



Ghazal Khalighinejad



Anej Svete



Josef Valvoda



Ryan Cotterell

#### **Definition: Alphabet**

An alphabet is a non-empty finite set of elements called symbols.

#### Examples:

- { 0, 1 }{ 0, 1, # }
- the set of all Unicode characters
- { above, accept, admire, ..., zebra, zebras, ., ? }

### **Definition: String**

A **string** is finite sequence of symbols from an alphabet.

#### Examples:

- 01011
- 01011#01011
- the salamanders don't amuse my newt.

#### **Definition: Language**

A language or formal language is a (possibly infinite) set of strings.

#### Examples:

```
• \{ #, 0 # 0, 1 # 1, 0 0 # 0 0, 0 1 # 0 1, 1 0 # 1 0, 1 1 # 1 1, ... \}
= \{ w#w \mid w \in \{0, 1\}^* \}
```

```
    { the salamanders don't amuse my newt . ,
        our zebra doesn't applaud the unicorn . ,
        some newt doesn't comfort the ravens . ,
        ... }
```

#### Membership in a Language

```
\{ w#w \mid w \in \{0, 1\}^* \}
```

- 0 1 1 # 0 1 1
- X 011#010
- X #0##10

### Natural Language as a Formal Language

- ✓ your orangutans giggle
- X orangutans giggle your

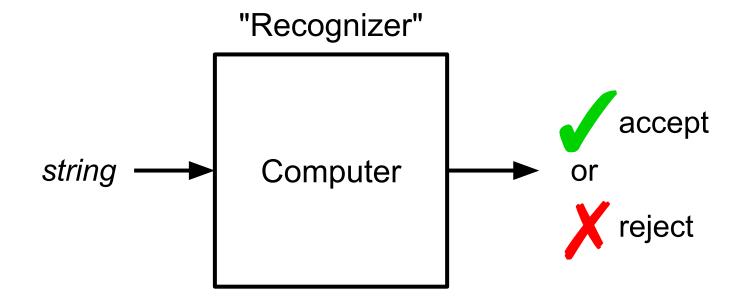
- ✓ your vulture comforts our peacocks by our newts
- × peacocks by comforts newts our your our vulture

- ✓ some newts that our zebra amuses wait
- X some newts that our zebra amuses waits

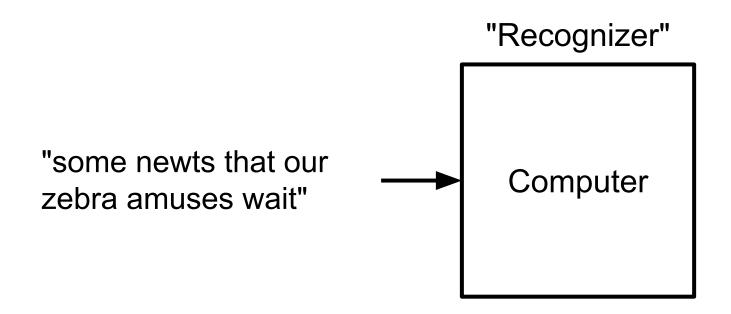
#### **Definition: Recognizer**

A **recognizer** is any computing device that reads a string as input and produces a decision to **accept** or **reject** it as output.

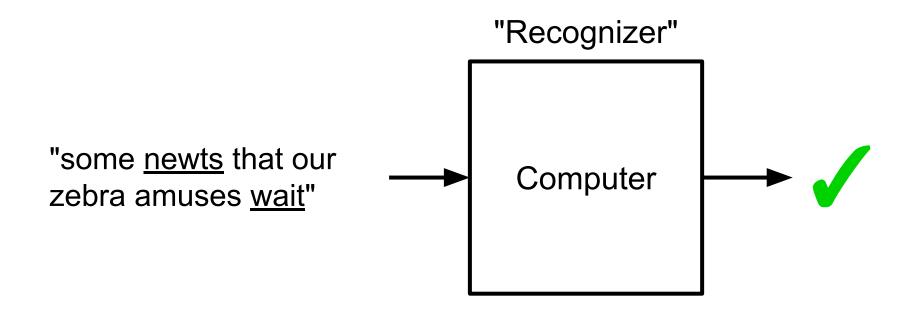
A recognizer recognizes the language of strings that it accepts.



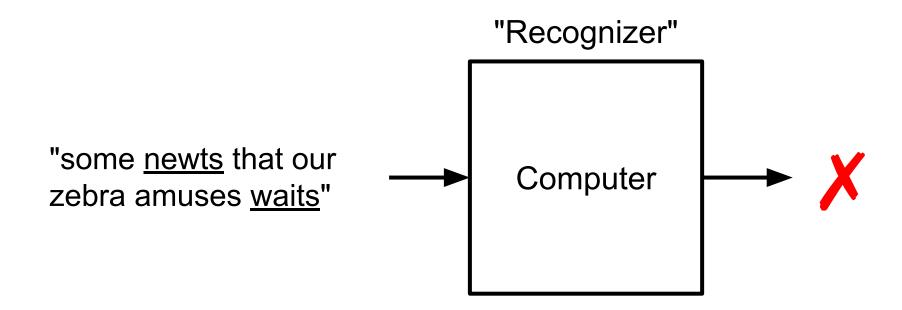
## Recognizers



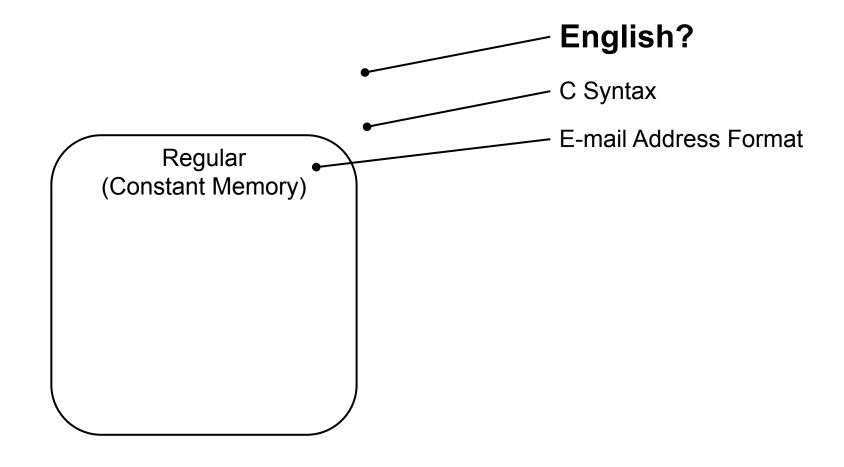
## Recognizers

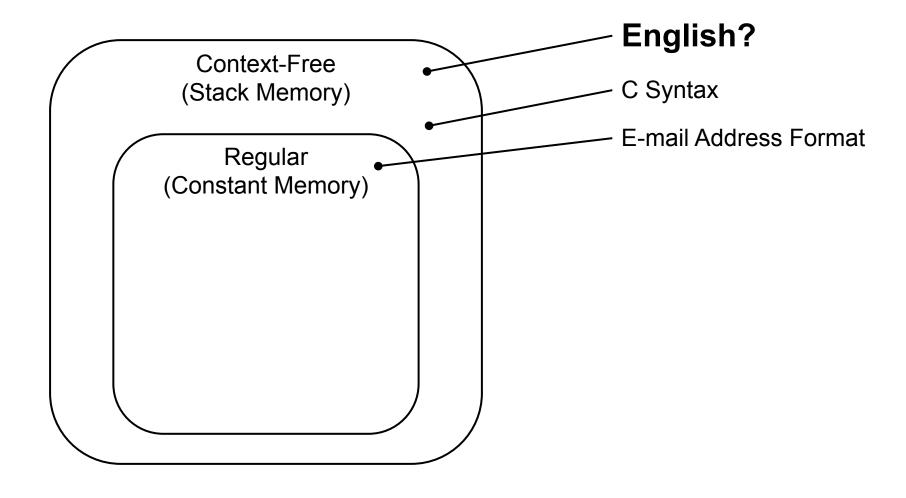


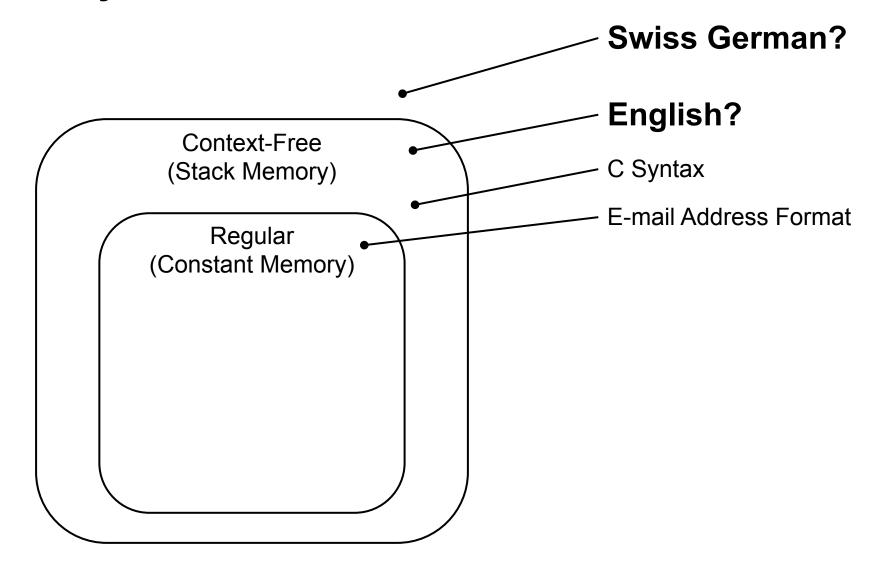
#### Recognizers

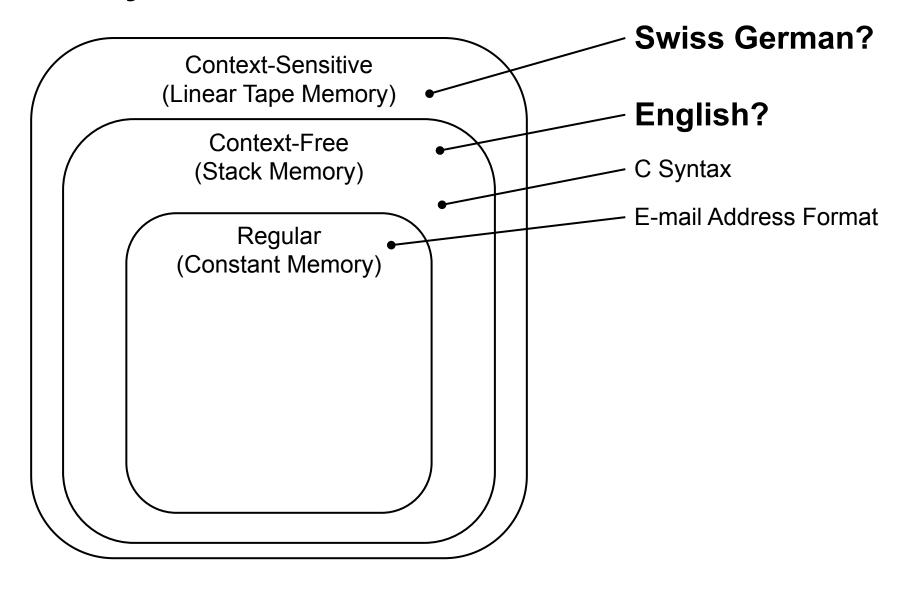


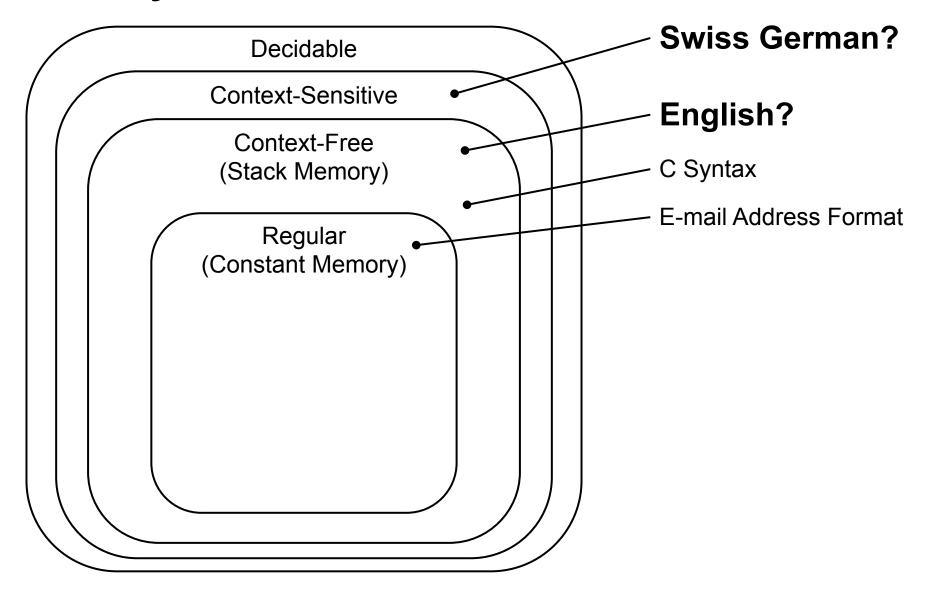
Regular (Constant Memory)

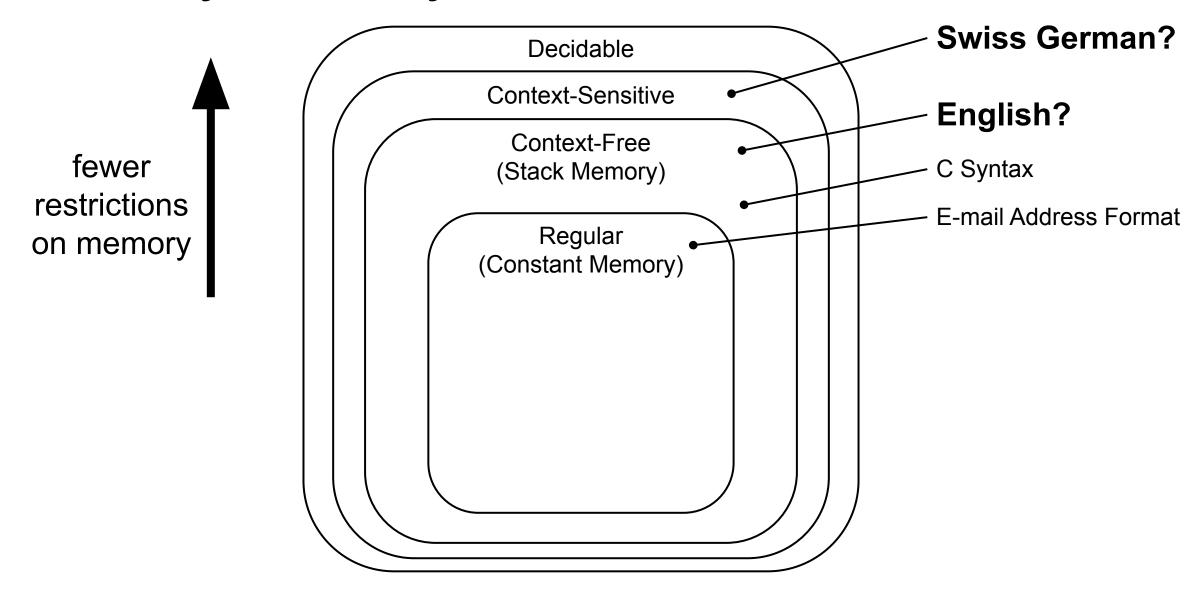






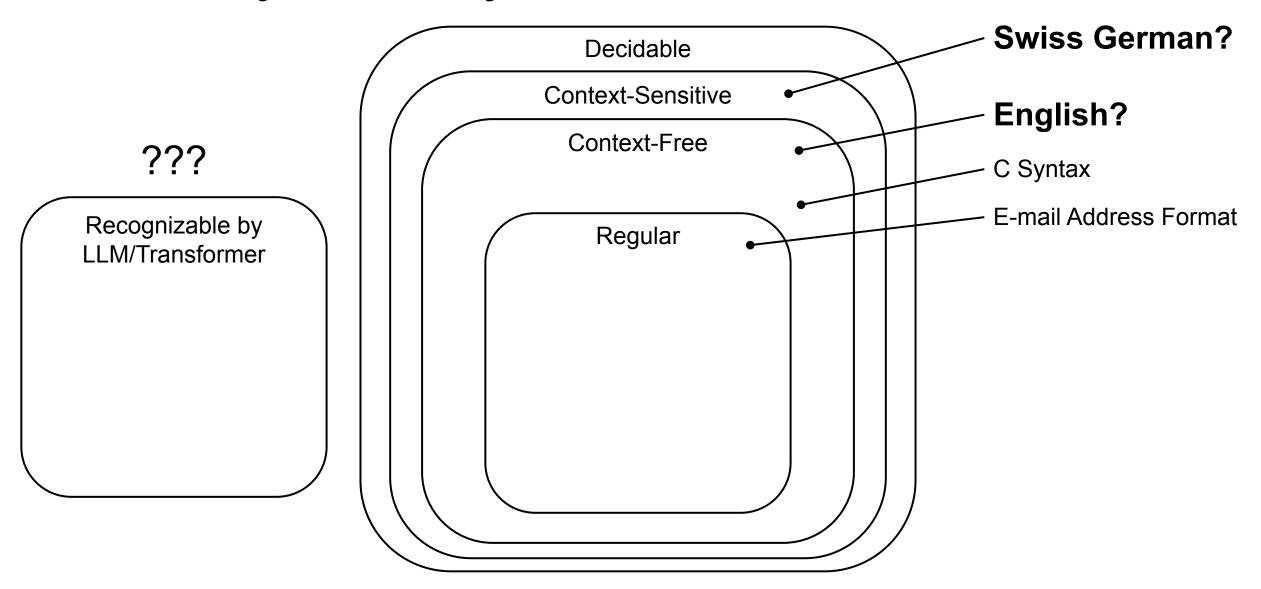




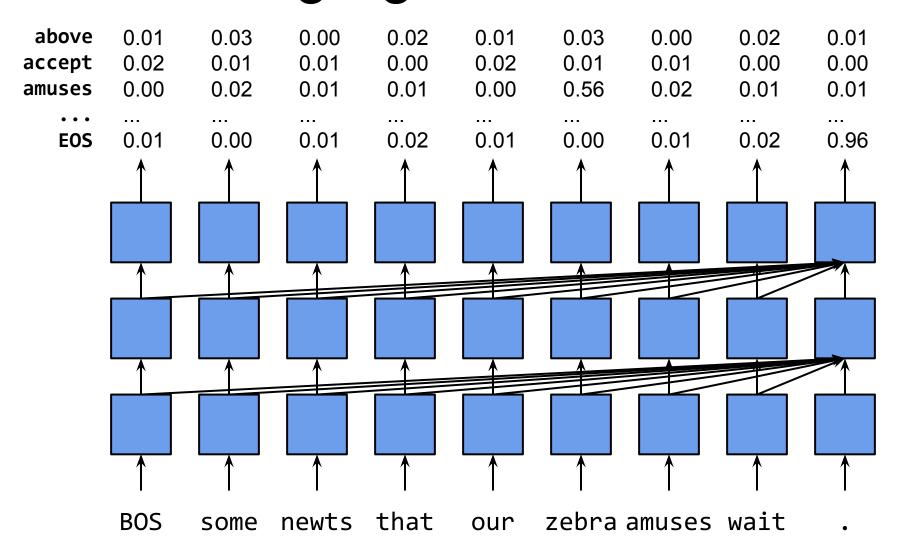


#### **Definition: Class**

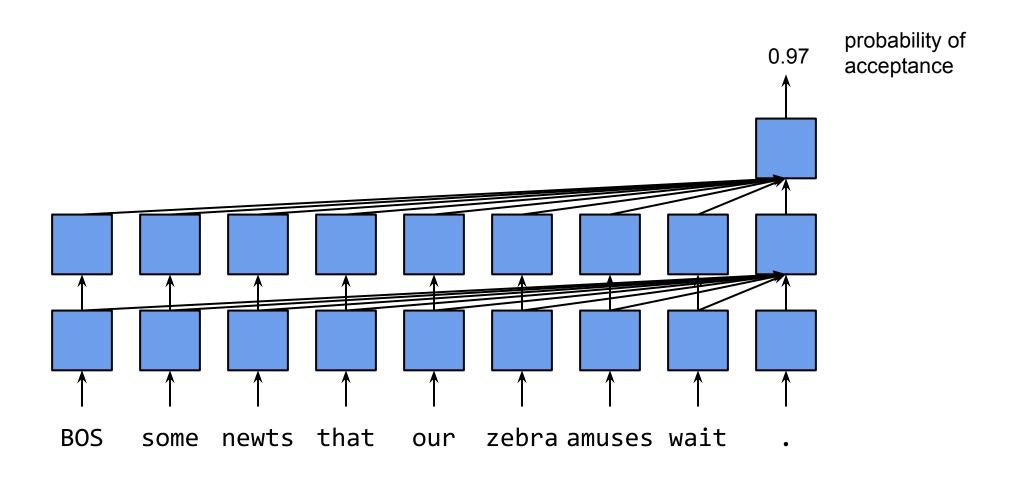
A class or language class is a (possibly infinite) set of languages.

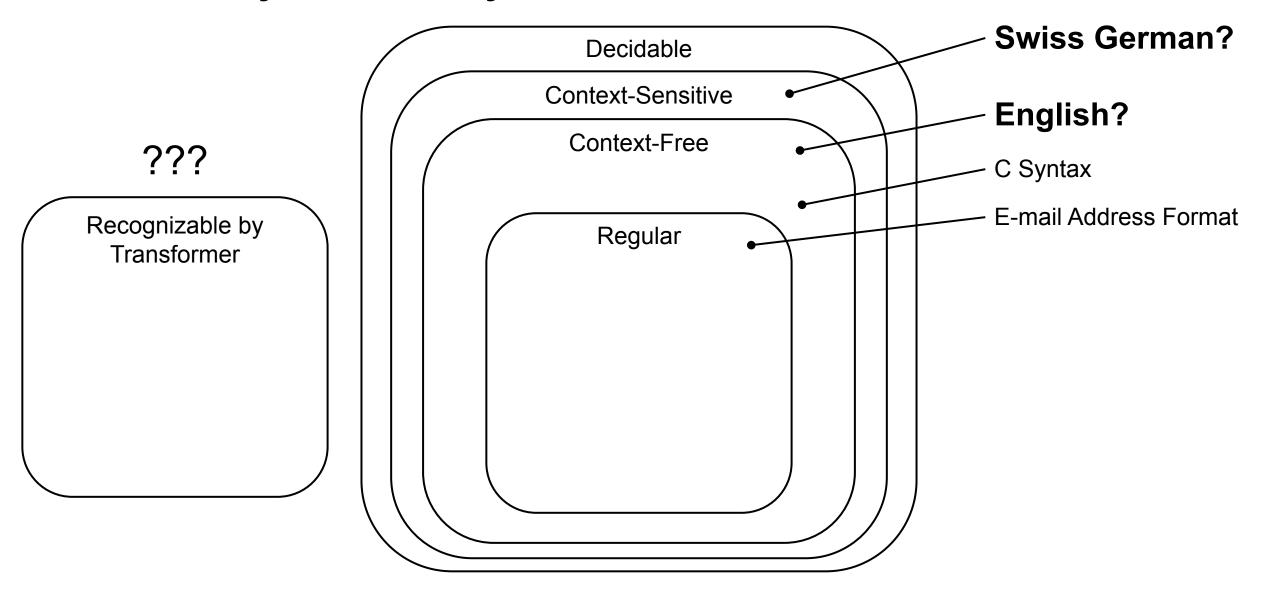


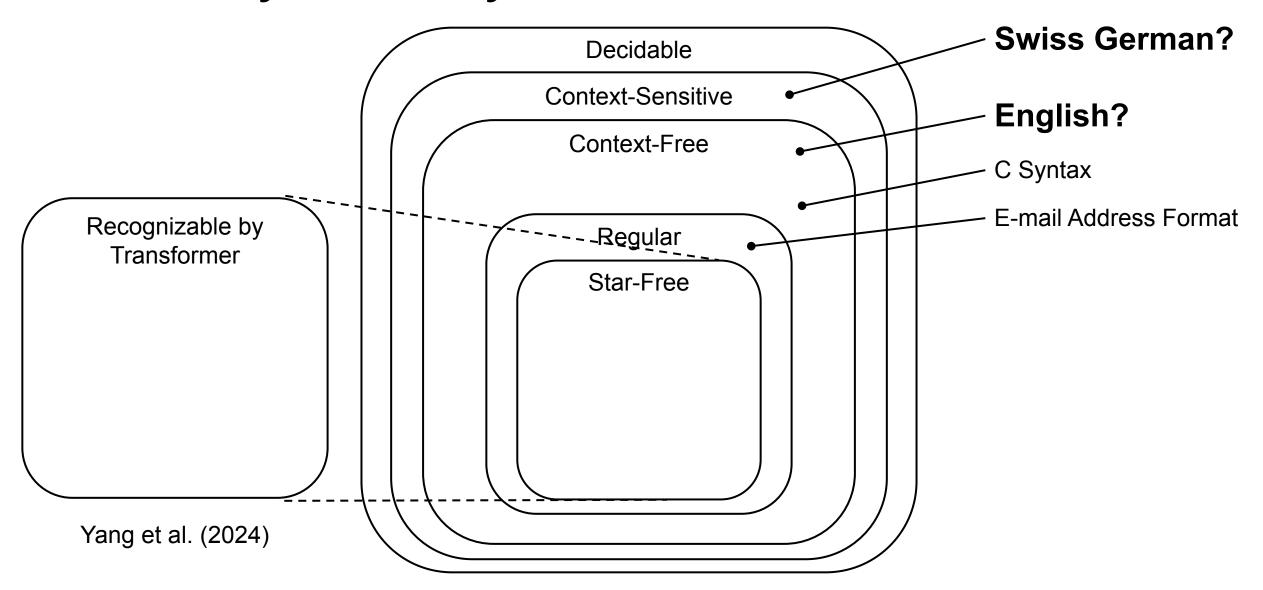
#### Transformer Language Model

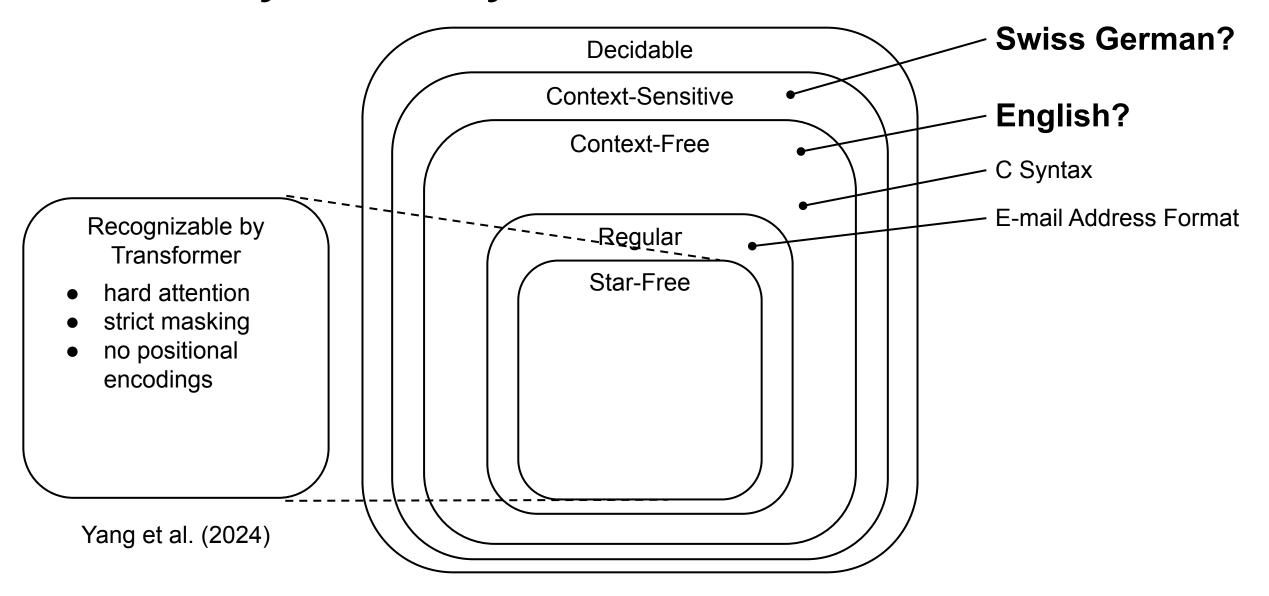


### **Transformer Recognizer**

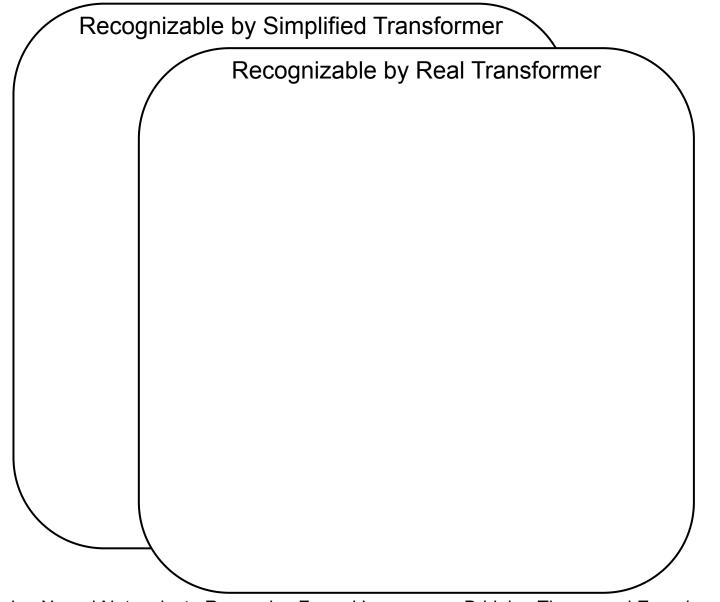




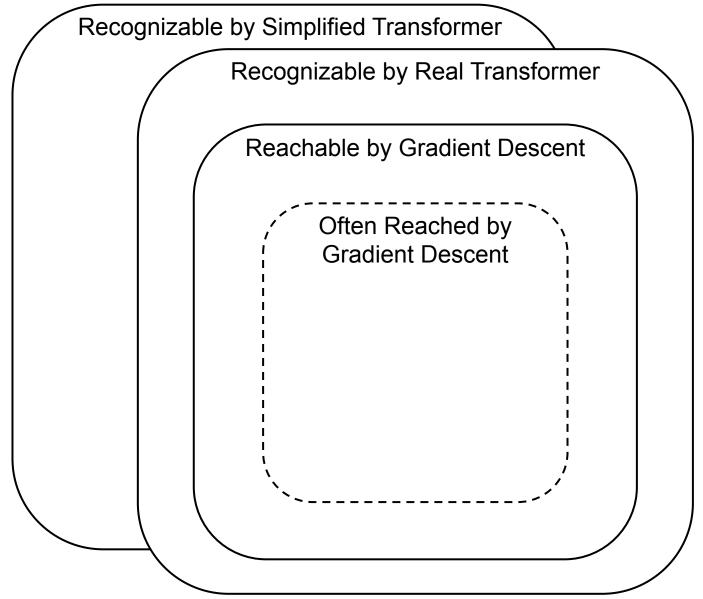




## **Transformer Expressivity**



#### **Transformer Expressivity**



#### **Purpose of This Work**

# How Do We Test Hypotheses about Neural Networks as Recognizers?

#### **Our Solution**

# Train Neural Networks as Recognizers (Binary Classifiers of Strings)

#### **Basic Approach**

- Get a dataset of strings labeled accept (1) or reject (0)
- The network outputs a probability of acceptance
- Use an objective function that maximizes the accept probability for accept and minimizes it for reject

#### Disconnect between Experiments and Theory

- Generating negative examples (strings labeled reject) is hard
- Positive-only tasks
  - Language modeling
  - Sequence-to-sequence transduction
- Can't convert a neural LM to a recognizer
  - How do you distinguish between probability of EOS that is just close to 0 or should be exactly 0?

Published as a conference paper at ICLR 2023

#### NEURAL NETWORKS AND THE CHOMSKY HIERARCHY

Grégoire Delétang\*1 Anian Ruoss\*1 Jordi Grau-Moya1 Tim Genewein1 Li Kevin Wenliang1

Elliot Catt<sup>1</sup> Chris Cundy<sup>†2</sup> Marcus Hutter<sup>1</sup> Shane Legg<sup>1</sup> Joel Veness<sup>1</sup> Pedro A. Ortega<sup>†</sup>

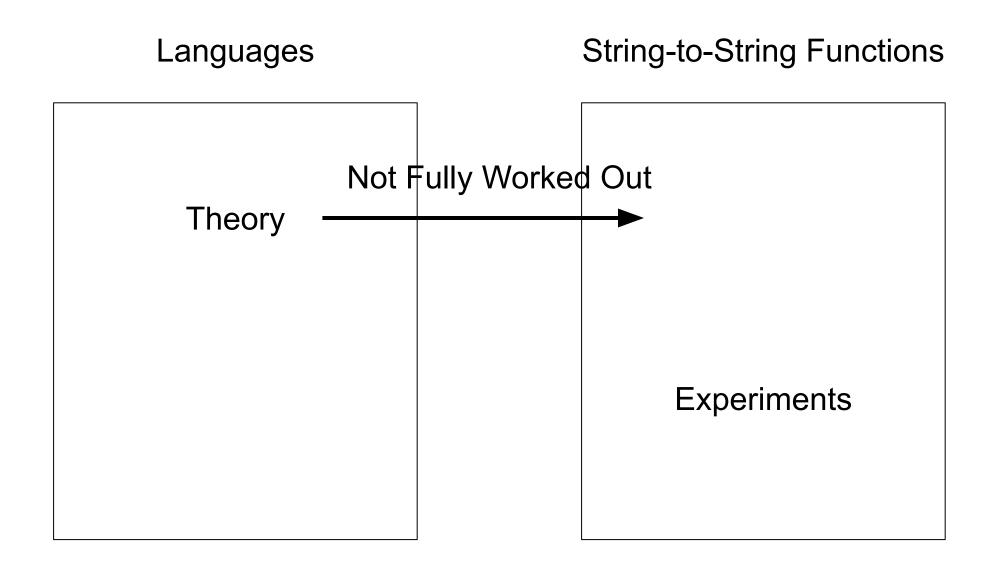
#### ABSTRACT

Reliable generalization lies at the heart of safe ML and AI. However, understanding when and how neural networks generalize remains one of the most important unsolved problems in the field. In this work, we conduct an extensive empirical study (20 910 models, 15 tasks) to investigate whether insights from the theory of computation can predict the limits of neural network generalization in practice. We demonstrate that grouping tasks according to the Chomsky hierarchy allows us to forecast whether certain architectures will be able to generalize to out-of-distribution inputs. This includes negative results where even extensive amounts of data and training time never lead to any non-trivial generalization, despite models having sufficient capacity to fit the training data perfectly. Our results

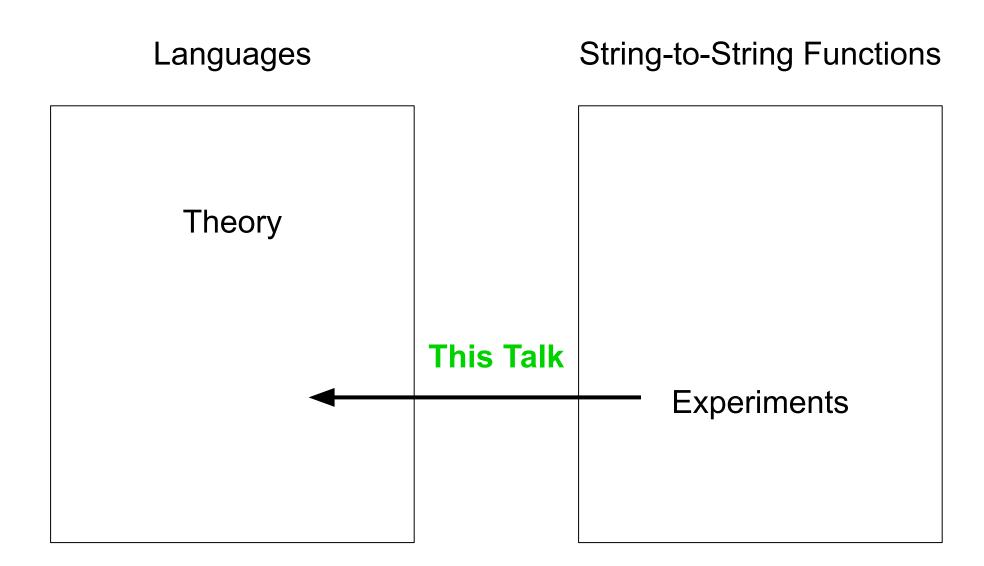
### **Bridging the Gap**

String-to-String Functions Languages Theory **Experiments** 

### **Bridging the Gap**



# **Bridging the Gap**



#### **Research Questions**

- How powerful is each NN architecture?
- Do our results change if we do recognition instead of string-to-string?
- What training objectives work best?

# Languages

### Languages

Table 1: Formal languages tested in this paper and included in FLaRe. For each language, we show the language class that it belongs to: regular (R), deterministic context-free (DCF), context-free (CF), or context-sensitive (CS). Each language does not belong to the previous language classes. Let  $c_u(w)$  be the number of times substring u occurs in w, let  $w_{i\to a}$  be w with its  $i^{\text{th}}$  symbol replaced with a, and let  $\langle x \rangle$  be the little-endian binary encoding of  $x \in \mathbb{Z}_{\geq 0}$ . See App. E for details.

Bigger	Language
Cla	asses

Class	Language	Description	Example String
R	Even Pairs	$ \{ w \in \{\emptyset, 1\}^* \mid c_{\emptyset 1}(w) + c_{1\emptyset}(w) \text{ is even} \} $ $= \{ aua \mid a \in \{\emptyset, 1\}, u \in \{\emptyset, 1\}^* \} \cup \{\varepsilon, \emptyset, 1\} $	010110
	Repeat 01	$\{(01)^n \mid n \ge 0\}$	010101
	Parity	$\{w \in \{0,1\}^* \mid c_1(w) \text{ is odd}\}$	11011001
	Cycle Navigation	A sequence of left (<), right (>), stay (=) moves on a 5-position cycle, then the final position (0-indexed).	>>=<>2
	Modular Arithmetic	Expression involving $\{+, -, \times\}$ and $\{0, \dots, 4\}$ , then the result mod 5. No operator precedence.	1-3×2=1
	$\operatorname{Dyck-}(2,3)$	Strings of balanced brackets with 2 bracket types and a maximum depth of 3.	()[([])()]
	First	$\{1w \mid w \in \{0,1\}^*\}$	100010
DCF	Majority	$\{w \in \{0,1\}^* \mid c_1(w) > c_0(w)\}$	101101
	Stack Manipulation	A stack from bottom to top, a sequence of push and pop operations, and the resulting stack from top to bottom.	011 POP =10
	Marked Reversal	$\{w \# w^R \mid w \in \{0,1\}^*\}$	001#100
CF	Unmarked Reversal	$\{ww^R \mid w \in \{0,1\}^*\}$	001100
CS	Marked Copy	$\{w\#w \mid w \in \{0,1\}^*\}$	001#001
	Missing Duplicate	$\{(ww)_{i\to_{-}} \mid w \in \{0,1\}^*, 1 \le i \le 2 w , (ww)_i = 1\}$	1_011101
	Odds First	$\{a_1b_1 \cdots a_nb_n a \# a_1 \cdots a_n a b_1 \cdots b_n \mid n \ge 0; a_i, b_i \in \{\emptyset, 1\}; a \in \{\emptyset, 1, \varepsilon\}\}$	01010=00011
	Binary Addition	$\{\langle x\rangle 0^i + \langle y\rangle 0^j = \langle x+y\rangle 0^k \mid x,y,i,j,k \in \mathbb{Z}_{\geq 0}\}$	110+01=10100
	Binary Multiplication	$\{\langle x\rangle 0^i \times \langle y\rangle 0^j = \langle xy\rangle 0^k \mid x, y, i, j, k \in \mathbb{Z}_{\geq 0}\}$	110×0100=011
	Compute Sqrt	$\{\langle x\rangle 0^i = \langle  \sqrt{x} \rangle 0^j \mid x, i, j \in \mathbb{Z}_{\geq 0}\}$	01010=1100
	Bucket Sort	Sequence of integers in $\{1, \dots, 5\}$ , then # and the sorted sequence.	45134#13445

### **Dataset Generation**

### Positive examples of w#w

- ✓ 01011#01011
- 110#110
- **/** #
- 100100#100100

### **Negative Examples of w#w**

- ✓ 01011#01011
- 110#110
- **/** #
- 100100#100100
- X #0#10##010#
- × 01#1#00
- X #11#0#0
  - 01#01

# **Negative Examples of w#w**

- ✓ 01011#01011
- 110#110
- **/** #
- 100100#100100
- X #0#10##010#
- X 01#1#00
- X #11#0#0

### **Negative Examples of w#w**

- ✓ 01011#01011
- X 01011#00011
- 110#110
- X 1100#110
- 100100#100100
- X 10010#100100

### **Sampling Datasets**

- Our method only needs two algorithms
  - a. Sample a positive string with length in the range  $[n_{min}, n_{max}]$
  - b. **Membership testing** -- is a string in the language?
- Balanced label distribution: 50% positive, 50% negative
- Not specific to any language class

# **Negative Sampling**

- Pick a length n uniformly from [n<sub>min</sub>, n<sub>max</sub>]
- Propose random strings, test whether they are in the language, reject positives
  - o 50% chance: Generate a random string of symbols of length n
    - Tend to be too easy
  - 50% chance: Generate a positive example of length n, apply K random edits to it
    - Adversarial

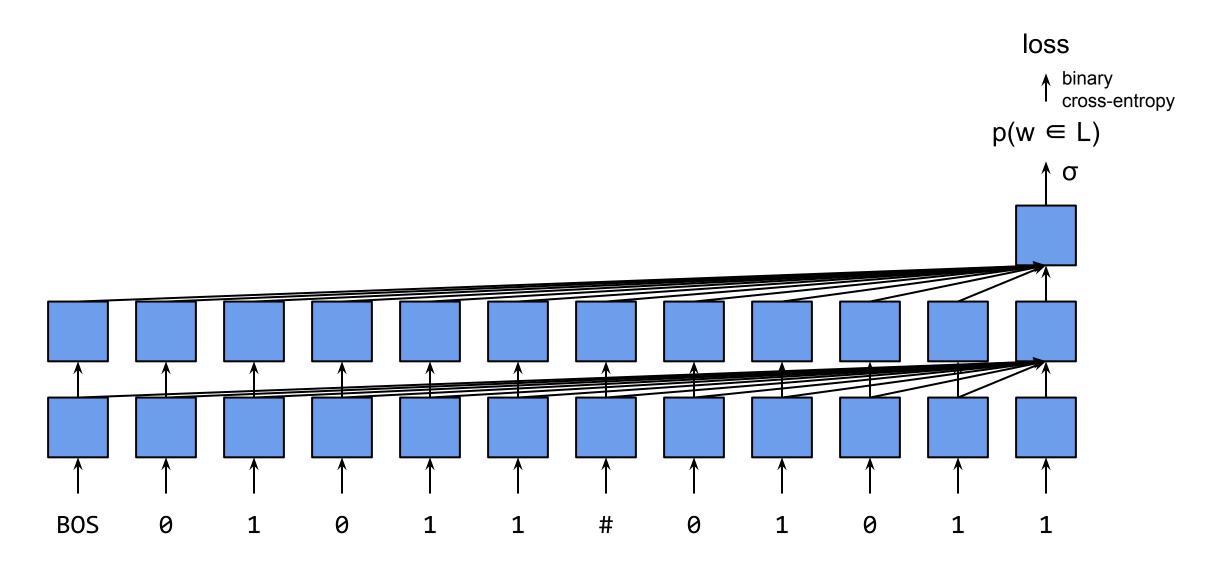
# Architectures and Training Objectives

#### **Three Architectures**

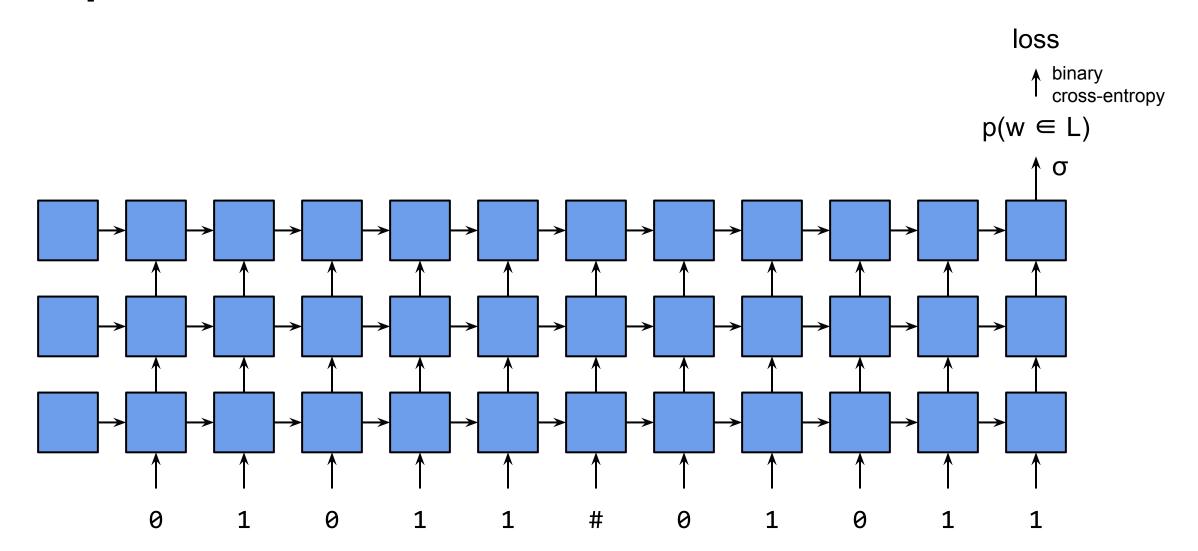
- Transformer
  - Causally masked (don't attend to later positions)
  - Sinusoidal positional encodings
  - Pre-norm instead of post-norm
- Simple RNN
  - tanh activation
  - learned initial state
- LSTM
  - decoupled input and forget gates
  - learned initial state

Every model has 5 layers and about 64k parameters.

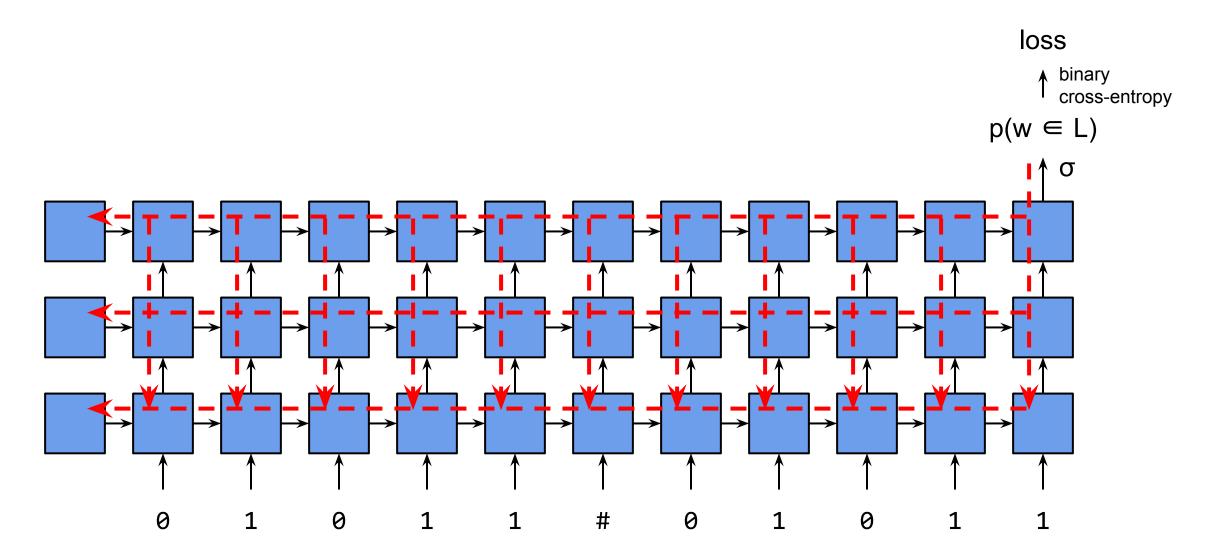
### **Transformer**



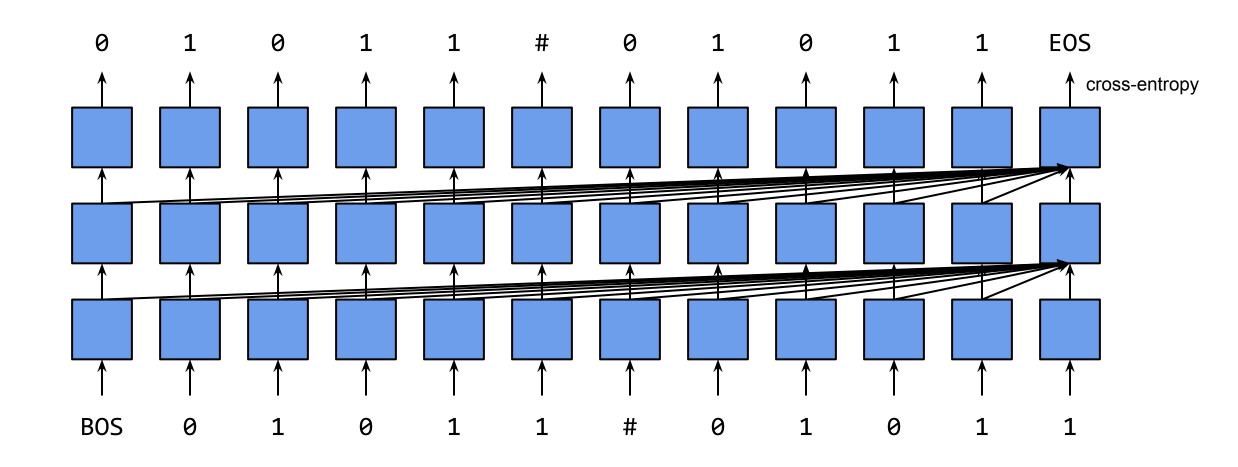
# Simple RNN and LSTM



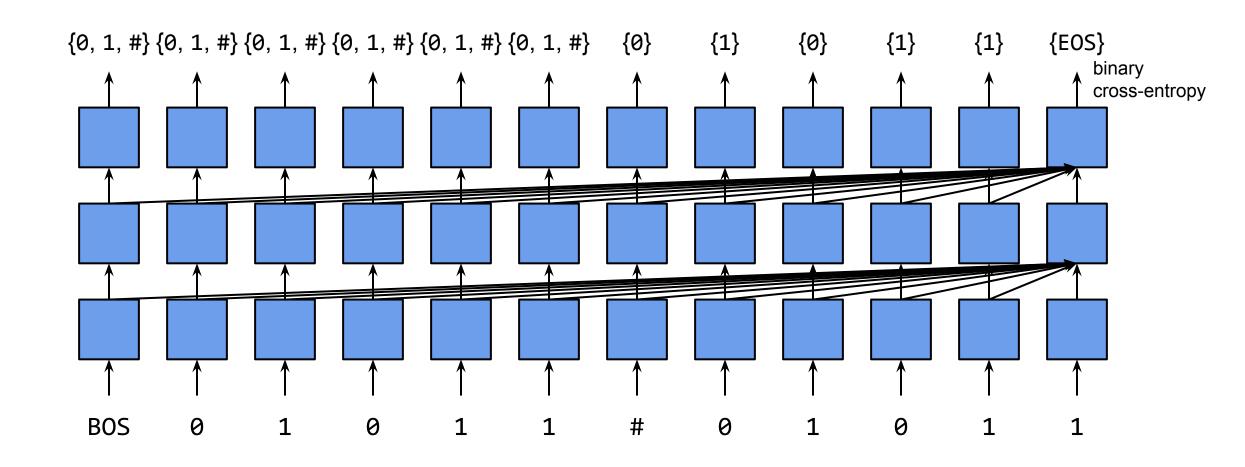
### **Gradient Problems**



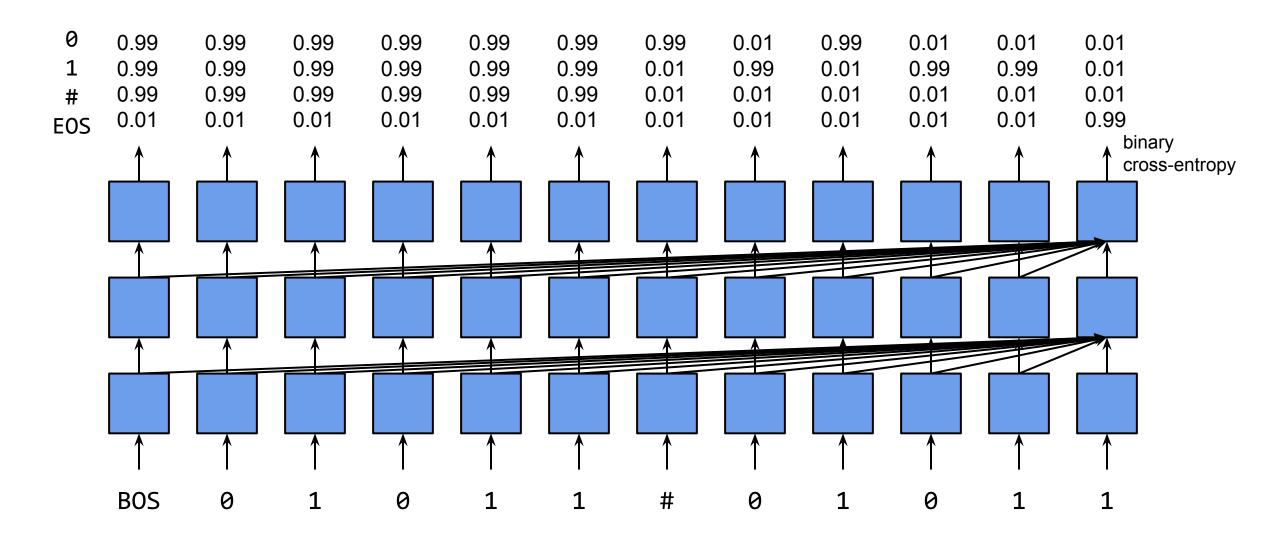
# Language Modeling



# **Next Symbol Prediction**



# **Next Symbol Prediction**



### **Multi-Task Learning**

Loss = Rec. Loss + 
$$\lambda_{LM}$$
 × LM Loss +  $\lambda_{NS}$  × NS Loss

- We always include the Rec. Loss term
- We test all 4 combinations of {with, without} × {LM Loss term, NS Loss term}

# Experiments

### **Dataset Splits**

- Training data: 10k examples with lengths in [0, 40]
- Test data: 5,010 examples with lengths in [0, 500]
  - Average of 10 examples per length
  - Emphasis on length generalization -- does it learn the algorithm?
- Two variations for validation data (next slides)

# **Testing Expressivity**

- Does a parameter setting exist at all?
- Long validation data set: 1k examples with lengths in [0, 80]
- Encourages model to length-generalize during training via
  - learning rate schedule
  - early stopping
- Report max test accuracy of 10 random seeds among all 4 loss functions



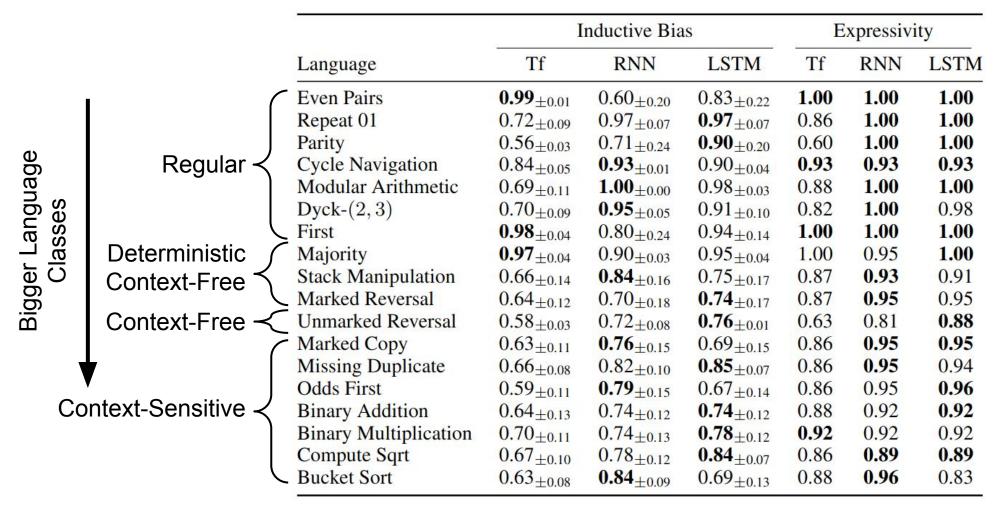
# **Testing Inductive Bias**

- What is the architecture's prior over languages?
- How does the network length-generalize without any hints during training?
- Short validation data: 1k examples with lengths in [0, 40]
- Report mean test accuracy of 10 random seeds of best loss function



#### Results

#### RNN and LSTM outperform Transformer; RNN is surprisingly good



### **Expressivity Results**

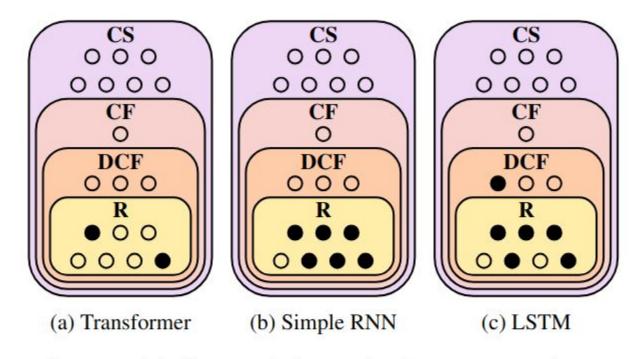
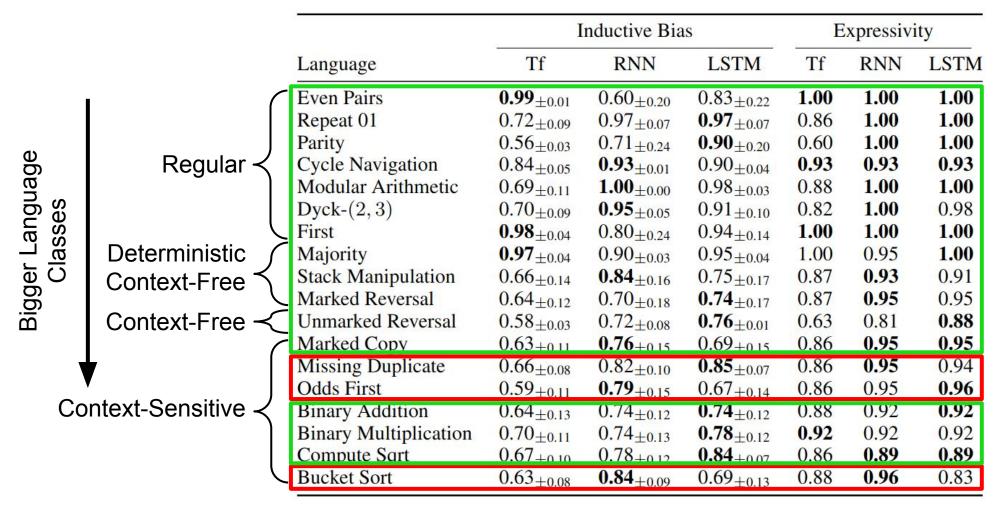


Figure 1: Summary of our empirical expressivity results. Dots represent languages, which are listed in Table 1. A filled dot means that the architecture exhibits perfect length generalization (see Table 2 under "Expressivity"). R = regular, DCF = deterministic context-free, CF = context-free, CS = context-sensitive. All architectures are limited to regular languages and the DCF language Majority. The transformer is strictly less expressive than the RNN/LSTM on the languages we tested.

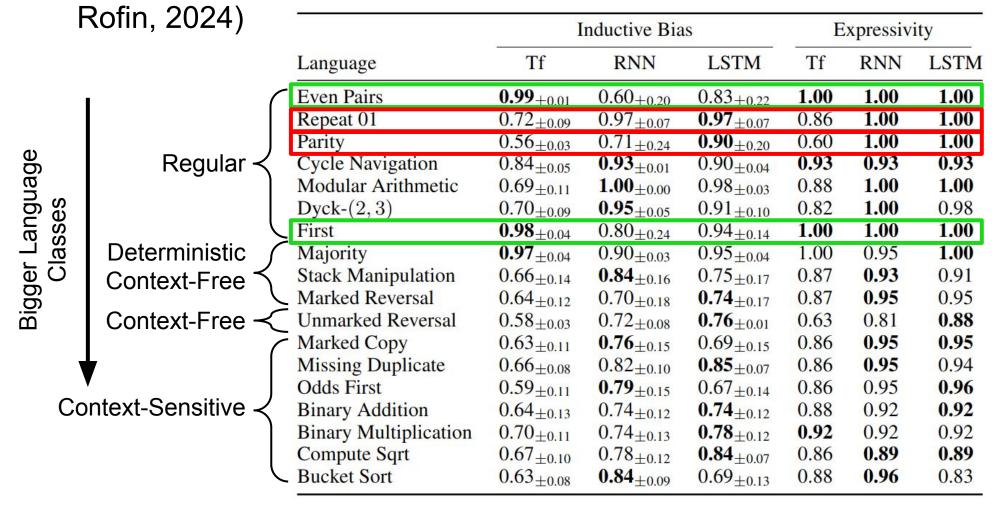
#### Results

#### Consistency in rankings between Inductive Bias and Expressivity



### Results

Transformers do well on low-sensitivity and badly on high-sensitivity (Hahn &



# **Low-Sensitivity**

**Even Pairs** (begins and ends with same bit)

```
<u>0</u> 1 0 1 1 <u>0</u>
```

<u>1</u>01<u>1</u>

First (begins with 1)

100101

<u>1</u>1011

<u>1</u> 0

# **High-Sensitivity**

Repeat 01 (0 1 repeated any number of times)

01010101

0101

Parity (odd number of 1's)

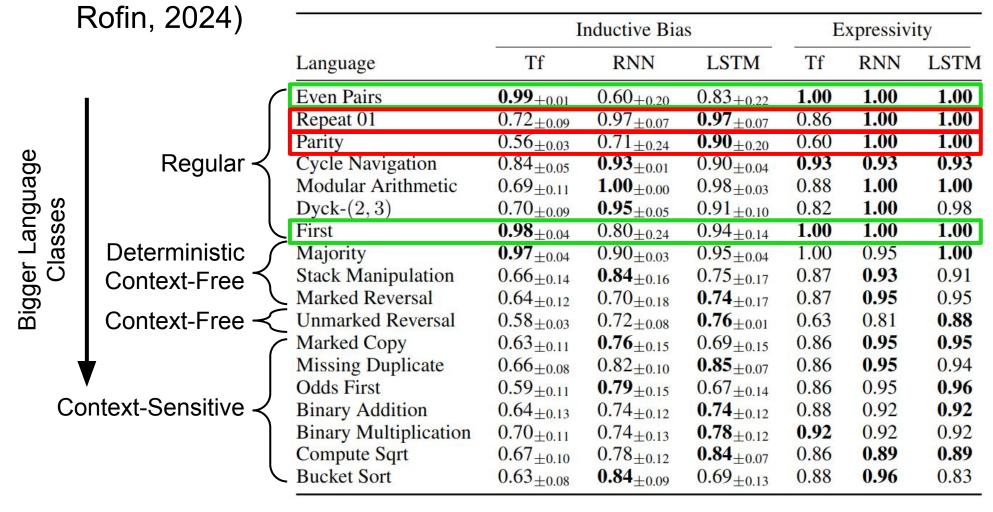
01011

1011

0001

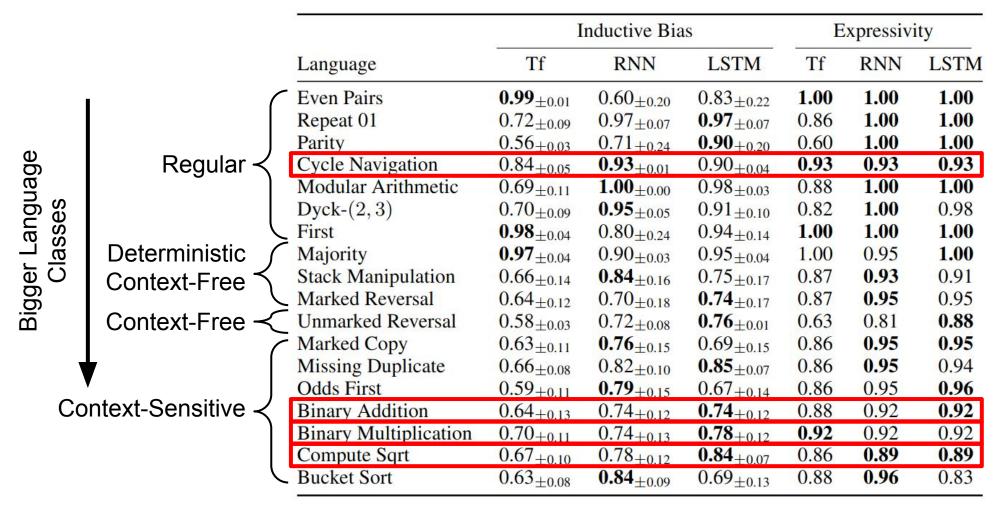
### Results

Transformers do well on low-sensitivity and badly on high-sensitivity (Hahn &

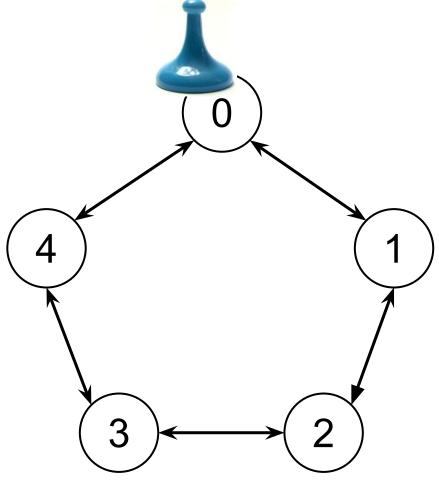


#### Results

#### Accuracy ceilings



**Cycle Navigation** 



### **Cycle Navigation**

```
√ >>>=><==>=>>><<>>>=<>=<=<=<=0
</p>
X = > > 2 1
><==>====0
× >===>><<===<<>3

√ <==><>=><<=<>2
X 1 > 4 0 2 2 1 0 3 4 0 2 > 1 4 > = 4 4 0 3 > < = 2 4 4 =
× >= 2 > 1 0 0 = < > 2 1 < 3 4 3 2 0 = 0 4 3 < 2 4 2 < 0 =
X <= < < > > < 0 < 2
✓ <> = < < 3</p>
```

### **Cycle Navigation**

- Increasing model size doesn't help
- Increasing training data size doesn't help
- Increasing adversarial negative examples doesn't help

# **Binary Arithmetic**

Small digits come first

#### **Addition**

$$0\ 0\ 1\ 0\ +\ 1\ 1\ 0\ =\ 1\ 1\ 1\ 0$$

#### Multiplication

$$0 1 0 0 * 0 0 1 = 0 0 0 1$$

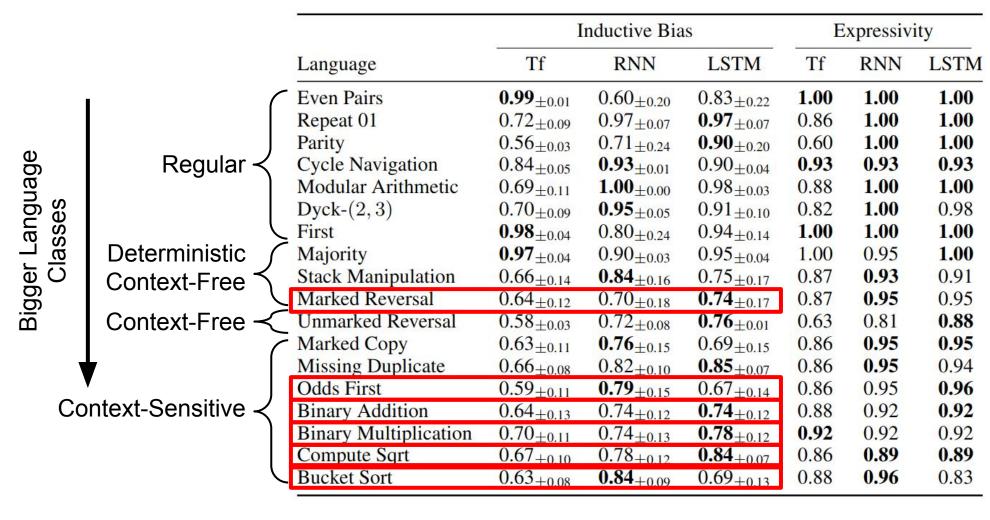
#### **Square Root**

# Recognition (Ours) vs. Seq-to-Seq (DeepMind)

- Transformer struggles on Parity
- Their RNN/LSTM solved Cycle Navigation, ours did not
- Our RNN performs better than theirs (slightly different architectures)

# Recognition (Ours) vs. Seq-to-Seq (DeepMind)

#### Different model rankings on many languages



### Performance vs. Edit Distance

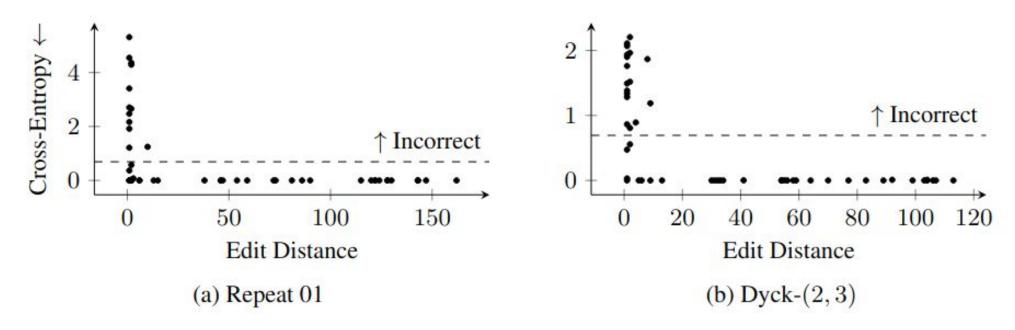
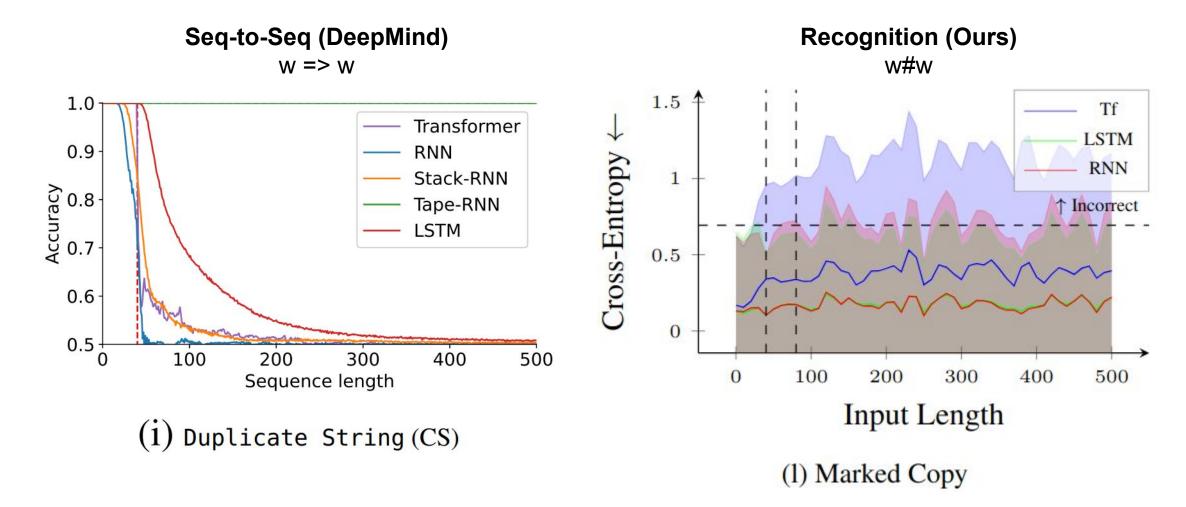
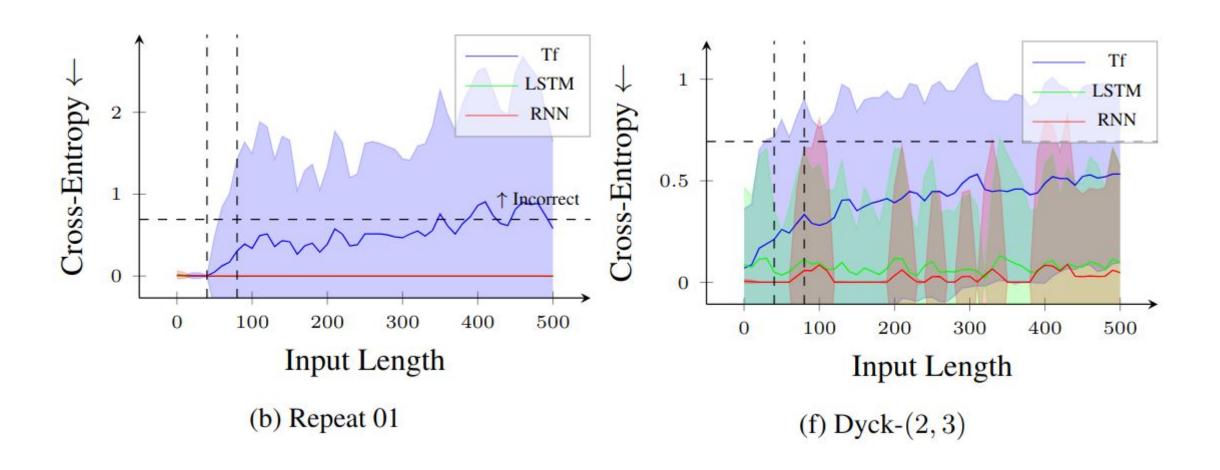


Figure 2: Recognition cross-entropy (lower is better) as a function of edit distance for the transformer model shown under "Expressivity" in Table 2, on a separate dataset of 50 negative examples in the length range [0,500]. The dashed lines show  $\log 2$ , the threshold for incorrect predictions. Despite being trained on a large proportion of negative examples with low edit distance, the transformer still struggles on examples that resemble positive examples.

### Performance vs. Length



### Performance vs. Length



### Which Loss Function Is Best?

Table 3: The best loss functions, corresponding to the accuracy scores reported in Table 2. "R" = recognition; "LM" = language modeling; "NS" = next symbol prediction. No single loss function consistently results in the best performance; the most frequent winner is just R.

	Inductive Bias			Expressivity		
Language	Tf	RNN	LSTM	Tf	RNN	LSTM
Even Pairs	R	R+LM+NS	R	R	R+NS	R
Repeat 01	R	R+NS	R	R	R	R
Parity	R+NS	R+NS	R+NS	R+NS	R+LM	R
Cycle Navigation	R+LM+NS	R	R	R	R	R
Modular Arithmetic	R	R	R	R+NS	R	R
Dyck-(2, 3)	R+LM	R+LM+NS	R	R+NS	R+NS	R+NS
First	R+NS	R+LM	R+LM	R	R	R
Majority	R+LM	R+NS	R+NS	R+LM	R+LM+NS	R+LM+NS
Stack Manipulation	R	R+NS	R	R+LM+NS	R	R+LM
Marked Reversal	R+NS	R+NS	R	R+LM	R+LM	R+LM
Unmarked Reversal	R	R+NS	R+NS	R	R+NS	R+NS
Marked Copy	R+NS	R+LM	R	R+NS	R	R
Missing Duplicate	R+LM+NS	R	R	R+LM+NS	R+LM+NS	R+LM
Odds First	R	R	R	R+LM+NS	R	R+LM+NS
Binary Addition	R	R+NS	R	R+LM	R+NS	R+NS
Binary Multiplication	R+NS	R	R+NS	R+NS	R+LM	R+NS
Compute Sqrt	R	R	R	R	R	R
Bucket Sort	R	R+LM+NS	R	R+NS	R+NS	R+LM+NS

#### What Did We Learn?

- Transformers, simple RNNs, and LSTMs rank very low in the Chomsky hierarchy
  - If true, they cannot make grammaticality judgments with 100% accuracy
- Simple RNNs and LSTMs outperform the transformer
- Using just a simple recognition objective is usually effective; auxiliary training objectives help in isolated cases but not uniformly
- Consistency between inductive bias and expressivity
- Results do change from the DeepMind paper
  - e.g., nothing is able to solve Cycle Navigation

### **Thank You**

Questions?

#### References

- Alexandra Butoi, Ghazal Khalighinejad, Anej Svete, Josef Valvoda, Ryan Cotterell, Brian DuSell. Training Neural Networks as Recognizers of Formal Languages. In Proc. ICLR. Singapore, May 2025.
- Andy Yang, David Chiang, Dana Angluin. Masked Hard-Attention Transformers Recognize Exactly the Star-Free Languages. In Advances in NeurIPS. Vancouver, Canada, December 2024.
- Grégoire Delétang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt, Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, Pedro A Ortega. Neural Networks and the Chomsky Hierarchy. In Proc. ICLR 2023. Kigali, Rwanda, May 2023.
- Michael Hahn, Mark Rofin. Why are Sensitive Functions Hard for Transformers? In Proc. ACL. Bangkok, Thailand, Aug 2024.