

# A Minimalist Dataset for Systematic Generalization of Perception, Syntax, and Semantics

Qing Li, Siyuan Huang, Yining Hong, Yixin Zhu, Ying Nian Wu, Song-Chun Zhu

(a) main concepts

syntax

number

parenthesis

op1

op]

op2 op2 semantics

0..5..9

none

 $i+j \\ \max(0,i-j)$ 

 $i \times j$ 

 $\operatorname{ceil}(i \div j)$ 

Semantics

Extra.

 $\checkmark$ 

 $\checkmark$ 

Inter.

V

perception

0.5.9

Syntax

Extra.

1

Inter.

 $\checkmark$ 

V



#### Can you decipher these ancient Egyptian signs?

へ気 合い気 ee → 60  つ ご 通 due → 100  い つ 通 u d → 12  の 取 の d 声 取 d → 4	ouうo台⊿⊿ → 18 U∽⊿uう台台 → 16 U&Пеஊ⊿ஊе⊿ → 41 ணuவிடுகபு → 4
$figerade mean \to 4$	ஊபฏி£ி⇔£⊿ → 4
⊠ ஊ $_{\Box}$ o n e e → $26$	၀ပဂည်းခ $\rightarrow 17$

Three-level Concept Learning and Generalization

#### syntax Perception: seen image $\rightarrow$ unseen image 45 ( Syntax: perception short expr. $\rightarrow$ long expr. 1+1=22×3=6 7-2=5 Semantics: 6×9=54 31-7=24 Small value $\rightarrow$ large value 17+23=40 Generalize semantics

#### Our contributions:

- We present HINT, a minimal yet comprehensive benchmark for systematic generalization w.r.t. perception, syntax, and semantics.
- Current neural networks, including Transformer and LLMs, struggle on HINT and the gap to human performance is considerable.
- Simply increasing model and dataset size does NOT lead to better generalization in HINT.

Code & dataset: <u>https://liqing-ustc.github.io/HINT</u>

# HINT: <u>Handwritten arithmetic with Int</u>egers

Input: handwritten expression	Output: result	
2X5÷9	2	05
5X5+(3-0-2)	32	() + -
4×(3+9)-7-10-5	) 41	× ÷

# **Train and Evaluation**

 $D_{\text{train}} \subset \mathcal{T}_{\text{train}} = \{(x, y) : |x| \leq 10, \max(v) \leq 100\},\$ 

Generalization: Interpolation & Extrapolation	Perce
$D_{\text{test}} = I \cup SS \cup LS \cup SL \cup LL$ , where	ption
$\mathrm{I} \subset D_{\mathrm{train}},$	$\checkmark$
$ extsf{SS} \subset \mathcal{T}_{ extsf{train}} ackslash D_{ extsf{train}},$	$\checkmark$
$LS \subset \{(x,y):  x  > 10, \max(v) \leqslant 100\}$	$\checkmark$
$\mathrm{SL} \subset \{(x,y):  x  \leqslant 10, \max(v) > 100\}$	$\checkmark$
$\mathrm{LL} \subset \{(x,y):  x  > 10, \max(v) > 100\}$	$\checkmark$

# **Examples from HINT**

Train		$2\times 5 \div 9 2 $ (9-9)×(3-4)-1×(0+3-(6-(2-2÷2))) 0
		5X5+(9-0-2) 32 4×(3+9)-7-(D-5) 41
	I	1=41 1×(2=5)×(8=8=6)0 6=4+(0=(6+0=(6+0)×1))+(9+4)15
	SS	1+3:4 2 3X(7×1)+(8+4)+4×3 66 4+(D-(7+7+7+6))×4-0 4
	LS	3×(8×(8×1)+0÷9)192 5×(3:4×9)+(2-5)×(7×(6+5))135
		2X(3X(3÷6+6X(3×4×6÷(1×6)))+O÷≥) <b>438</b>
Test	SL	(6×5-0)÷((4+3+5)÷9)+(3-((2-(2+(3×7-8÷9)))/4-9)) 18
		$6-3 \div (9 \times (9 \div (4-(4-7)))) \div (1 \div (7 \times 2 \div 6 \div 8))$ 6
		$(7+3)/(6-6\times(0\times(6+1)))-(3\times(-6-4)(4-3))\times(9\times3)$ 2
	LL	(6+)×(+2:4+(+++-0+3)×8-(1+3×8))×((0+(2×8-0)/3)+(8+9))174
		(3+(8+(4-7)(++8))×(8+4-(4-(6+5)+6))))÷(7+5×1×0)+5 1
		7X(8÷(1X(7÷+))+(1+>))X10+9-5÷[8+4÷19×6)))+18-(9-8+3))620

# **Experimental Results**

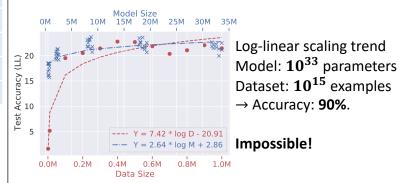
#### **Image Inputs**

Model	Variant	I	SS	LS	SL	LL	Avg.
GRU	w/o att	61.3±1.4	$53.3 \pm 1.7$	$30.5 \pm 1.2$	$9.2 \pm 0.2$	$11.9 \pm 0.5$	33.2±0.9
	w/ att	$66.7 \pm 2.0$	$58.7 \pm 2.2$	$33.1 \pm 2.7$	$9.4 \pm 0.3$	$12.8 \pm 1.0$	$35.9 \pm 1.6$
LSTM	w/o att	$80.0\pm5.7$	$76.2 \pm 7.4$	$55.7 \pm 8.2$	$10.9 \pm 0.6$	$19.8 \pm 2.6$	$48.6 \pm 4.9$
LSIM	w/ att	$83.9 \pm 0.9$	$79.7 \pm 0.8$	$62.0 \pm 2.5$	$11.2 \pm 0.1$	$21.0 \pm 0.8$	$51.5 \pm 1.0$
	vanilla	$20.9 \pm 0.4$	9.3±0.2	$5.7 \pm 0.3$	$1.5 \pm 0.3$	$2.9 \pm 0.5$	8.3±0.3
Transformer	rel.	$86.2 \pm 0.9$	$83.1 \pm 1.3$	$60.1 \pm 2.3$	$10.9 \pm 0.2$	$19.4 \pm 0.5$	$51.7 \pm 1.0$
	rel. uni.	88.4±1.3	86.0±1.3	62.5 <u>+</u> 4.1	$10.9\pm0.2$	$19.0 {\pm} 1.0$	$53.1\pm1.6$

#### Symbol Inputs

Model	Variant	I	SS	LS	SL	LL	Avg.
GRU	w/o att	$74.9 \pm 1.6$	$68.1 {\pm} 0.5$	$42.1 \pm 1.9$	$10.5 {\pm} 0.2$	$14.0\pm0.8$	41.3±0.6
GRU	w/ att	$76.2 \pm 0.6$	69.5 <u>±</u> 0.6	$42.8 \pm 1.5$	$10.5 \pm 0.2$	$15.1 \pm 1.2$	$42.5 \pm 0.7$
LSTM	w/o att	$84.3 \pm 5.2$	79.6±6.0	$63.7 \pm 6.1$	$11.7 \pm 0.3$	$22.1 \pm 1.4$	$52.3 \pm 3.8$
LSTW	w/ att	$92.9 \pm 1.4$	$90.9 \pm 1.1$	$74.9 \pm 1.5$	$12.1 \pm 0.2$	$24.3 \pm 0.3$	$58.9 \pm 0.7$
	vanilla	93.9±0.3	91.0±0.5	$33.2 \pm 1.2$	$11.5 \pm 0.1$	$11.5 \pm 0.7$	$47.4 \pm 0.4$
Transformer	rel.	$96.6 \pm 0.3$	$95.1 \pm 0.4$	$72.1 \pm 1.5$	$11.8 {\pm} 0.2$	$22.3 \pm 0.6$	$59.4 \pm 0.5$
	rel. uni.	98.0±0.3	96.8±0.6	78.2±2.9	$11.7 \pm 0.3$	$22.4 \pm 1.1$	61.5±0.9
GPT-3	0-shot	19.0	9.0	3.0	10.0	2.0	8.6
Gr 1-5	0-CoT	42.0	36.0	5.0	49.0	6.0	27.6

# Scaling laws w.r.t. model and dataset



#### **Few-shot Learning and Generalization**

