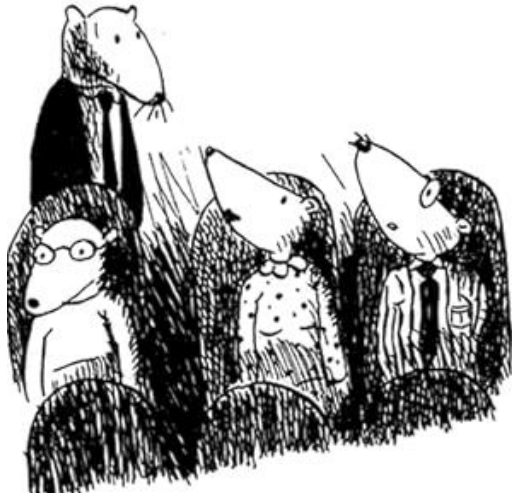


GROKKing

Ameli Alaeva

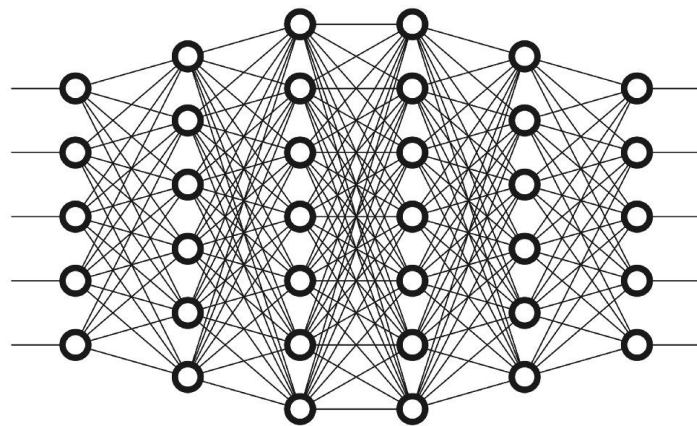




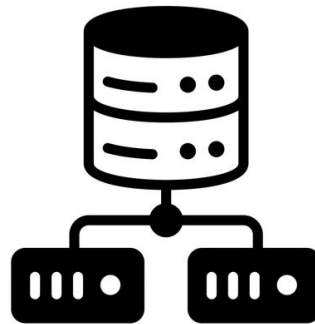
1



1



generalization



Dataset

PLAN

- reference
- 1st paper
 - observe dataset + model architecture
 - the goal
 - experiments + results

- 2d paper
- additional info
- conclusion

Q&A in the end





grok

/ɡrɒk/

verb

INFORMAL • US

gerund or present participle: **grokking**

understand (something) intuitively or by empathy.

"corporate leaders seemed to grok this concept fairly quickly"

GROKking: GENERALIZATION BEYOND OVERFITTING ON SMALL ALGORITHMIC DATASETS

Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin
OpenAI

Vedant Misra*
Google

arXiv:2201.02177v1 [cs.LG] 6 Jan 2022

All of our experiments used a small transformer trained on datasets of equations of the form $a \circ b = c$, where each of “ a ”, “ \circ ”, “ b ”, “ $=$ ”, and “ c ” is a separate token. Details of the operations studied, the architecture, training hyperparameters and tokenization can be found in Appendix A.1.

A.1.1 BINARY OPERATIONS

The following are the binary operations that we have tried (for a prime number $p = 97$):

$$x \circ y = x + y \pmod{p} \text{ for } 0 \leq x, y < p$$

$$x \circ y = x - y \pmod{p} \text{ for } 0 \leq x, y < p$$

$$x \circ y = x/y \pmod{p} \text{ for } 0 \leq x < p, 0 < y < p$$

$$x \circ y = [x/y \pmod{p} \text{ if } y \text{ is odd, otherwise } x - y \pmod{p}] \text{ for } 0 \leq x, y < p$$

$$x \circ y = x^2 + y^2 \pmod{p} \text{ for } 0 \leq x, y < p$$

$$x \circ y = x^2 + xy + y^2 \pmod{p} \text{ for } 0 \leq x, y < p$$

$$x \circ y = x^2 + xy + y^2 + x \pmod{p} \text{ for } 0 \leq x, y < p$$

$$x \circ y = x^3 + xy \pmod{p} \text{ for } 0 \leq x, y < p$$

$$x \circ y = x^3 + xy^2 + y \pmod{p} \text{ for } 0 \leq x, y < p$$

$$x \circ y = x \cdot y \text{ for } x, y \in S_5$$

$$x \circ y = x \cdot y \cdot x^{-1} \text{ for } x, y \in S_5$$

$$x \circ y = x \cdot y \cdot x \text{ for } x, y \in S_5$$

For each binary operation we constructed a dataset of equations of the form $\langle x \rangle \langle op \rangle \langle y \rangle \langle = \rangle \langle x \circ y \rangle$, where $\langle a \rangle$ stands for the token corresponding to element a .

★	a	b	c	d	e
a	a	d	?	c	d
b	c	d	d	a	c
c	?	e	d	b	d
d	a	?	?	b	c
e	b	b	c	?	a



★	a	b	c	d	e
a	a	d	?	c	d
b	c	d	d	a	c
c	?	e	d	b	d
d	a	?	?	b	c
e	b	b	c	?	a

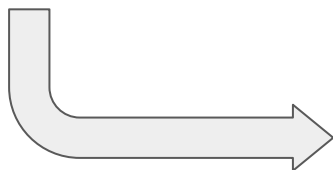


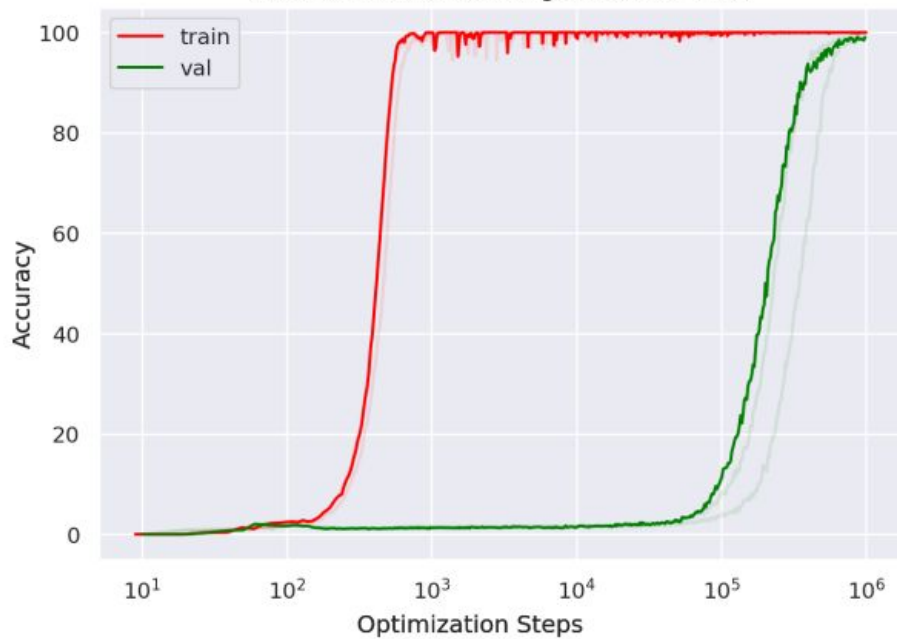
Figure 5: One of the binary operation tables presented to the networks that the network can perfectly fill in. Each symbol is represented as a letter in English, Hebrew, or Greek alphabet for reader's convenience. We invite the reader to guess which operation is represented here.

“a”, “b”, “c”, “=” and “o.”

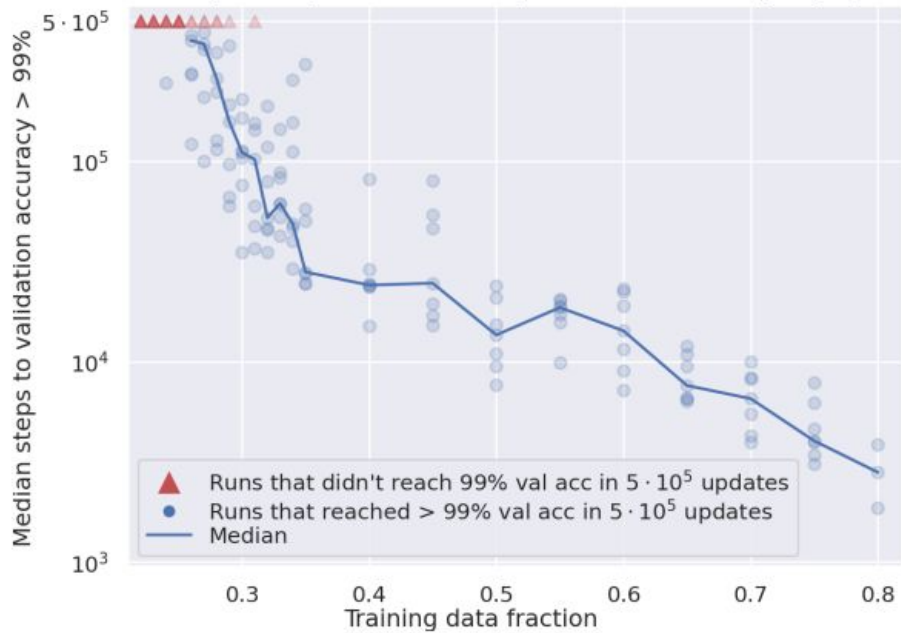
- a standard decoder-only transformer Vaswani et al. (2017)
- causal attention masking
- transformer with 2 layers,
- width 128
- 4 attention heads

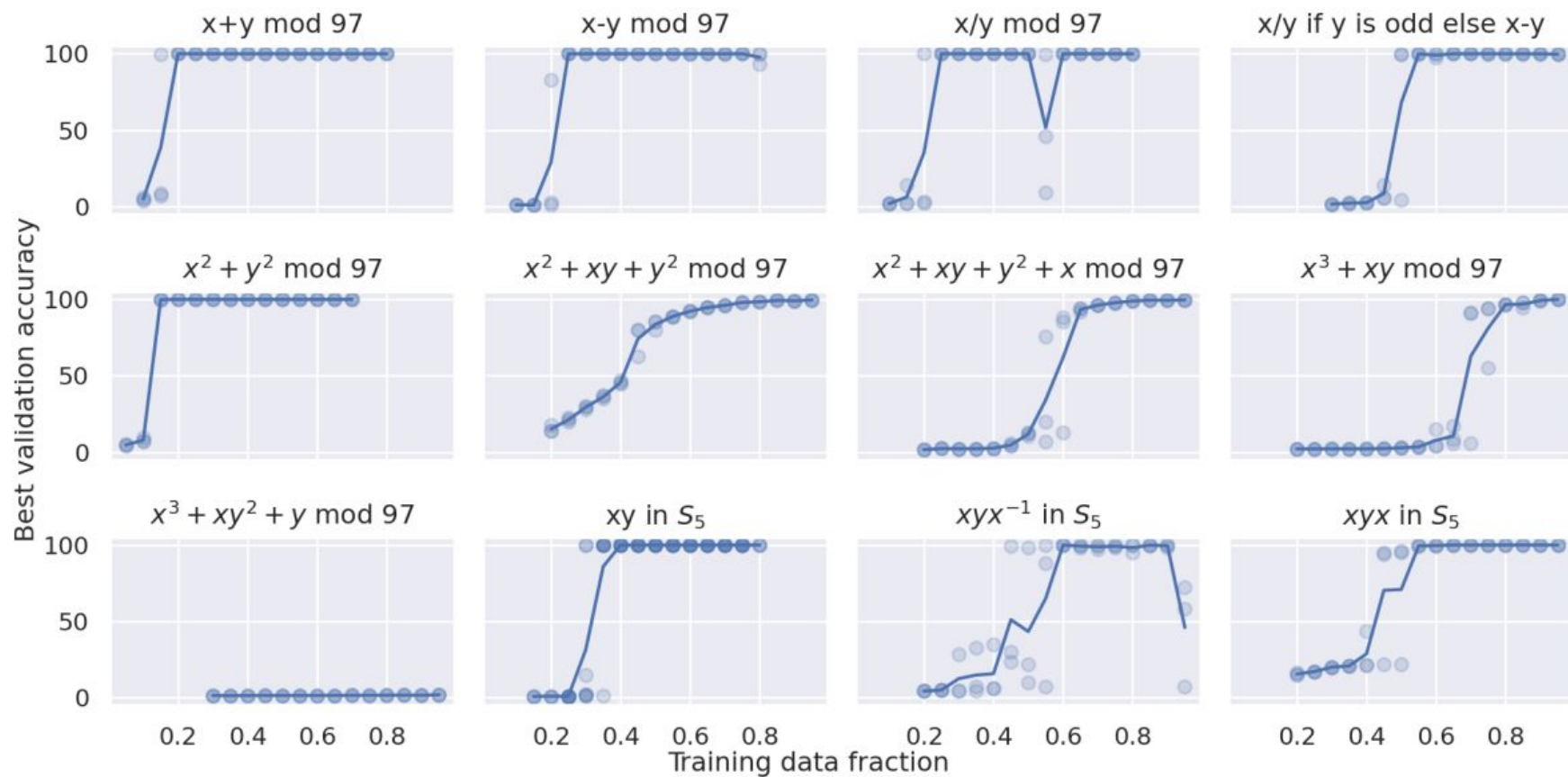
total of about $4 \cdot 10^5$ non-embedding parameters

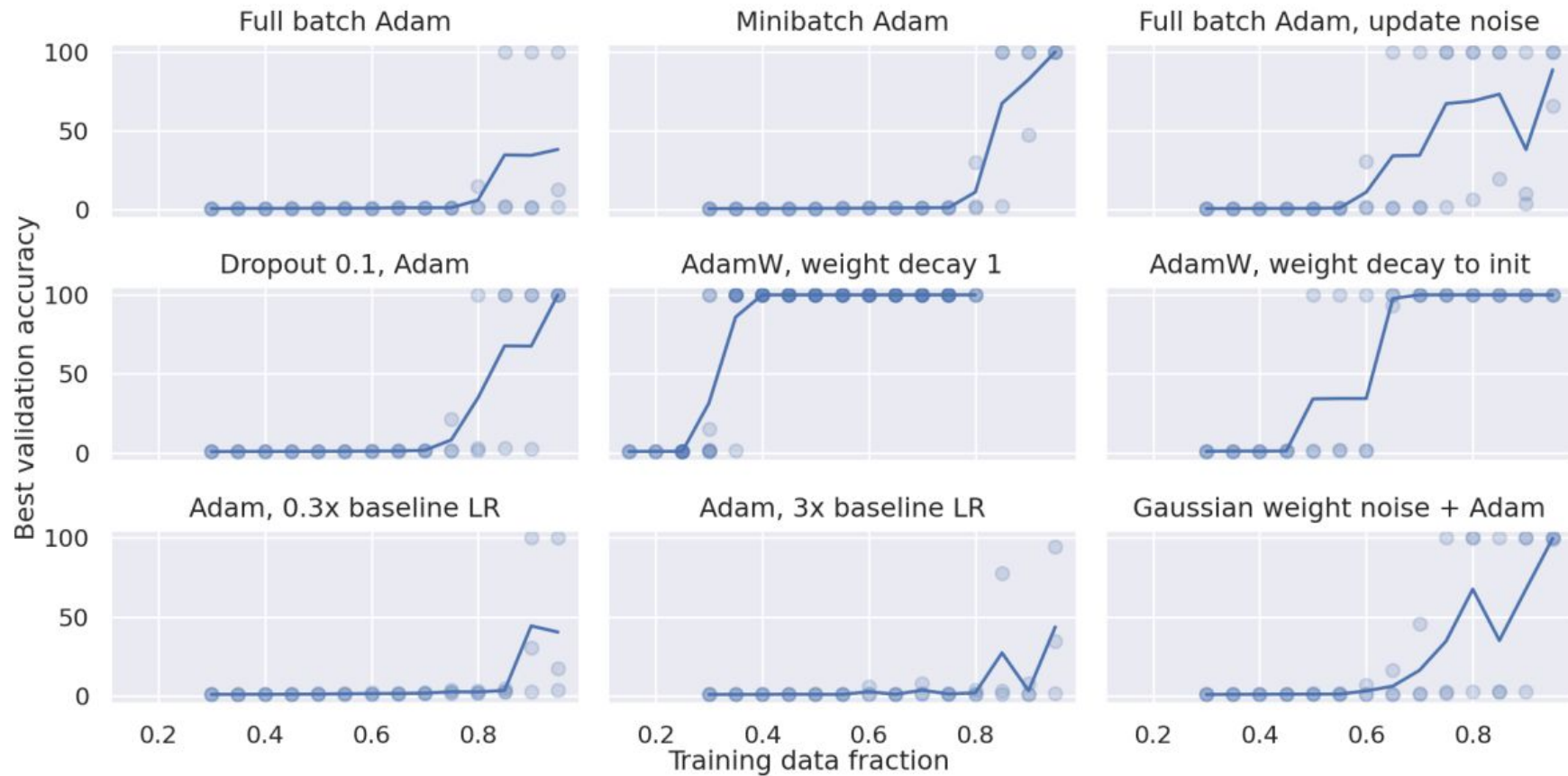
Modular Division (training on 50% of data)



Steps until generalization for product in abstract group S_5







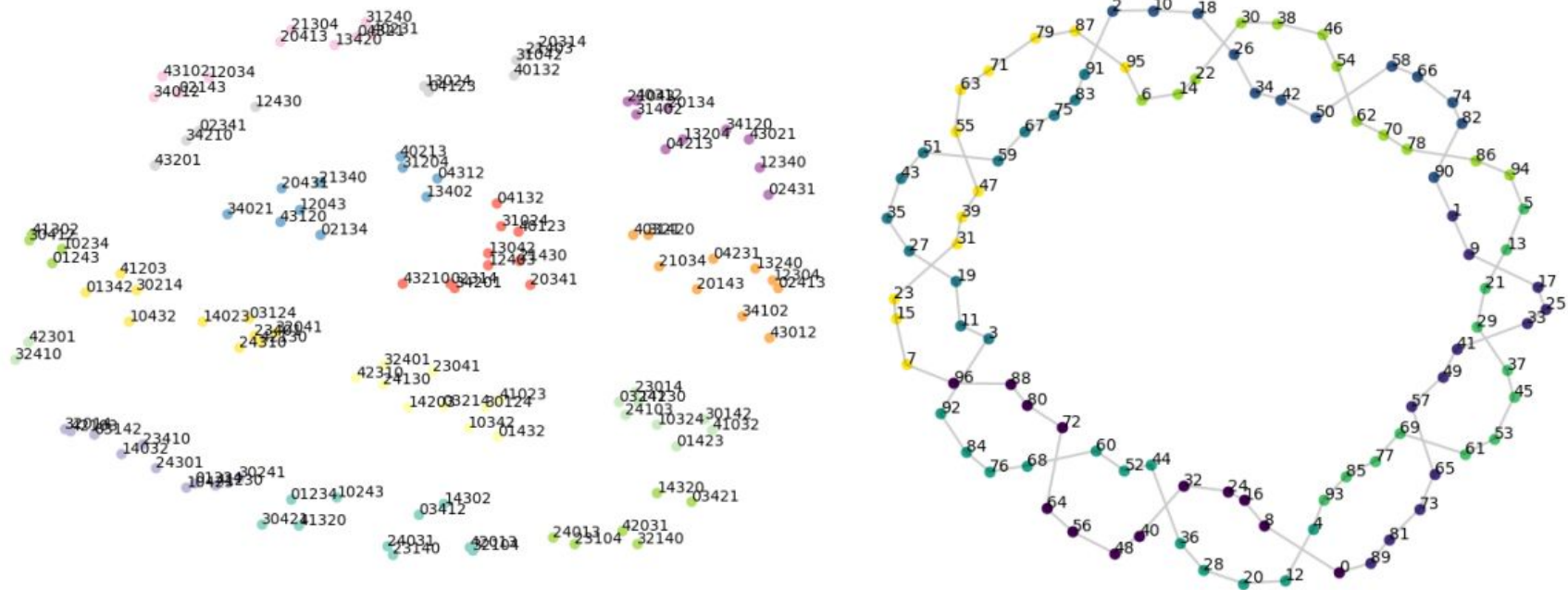


Figure 3: **Left.** t-SNE projection of the output layer weights from a network trained on S_5 . We see clusters of permutations, and each cluster is a coset of the subgroup $\langle (0, 3)(1, 4), (1, 2)(3, 4) \rangle$ or one of its conjugates. **Right.** t-SNE projection of the output layer weights from a network trained on modular addition. The lines show the result of adding 8 to each element. The colors show the residue of each element modulo 8.

PROGRESS MEASURES FOR GROKING VIA MECHANISTIC INTERPRETABILITY

Neel Nanda^{*,†} **Lawrence Chan[‡]** **Tom Lieberum[†]** **Jess Smith[†]** **Jacob Steinhardt[‡]**

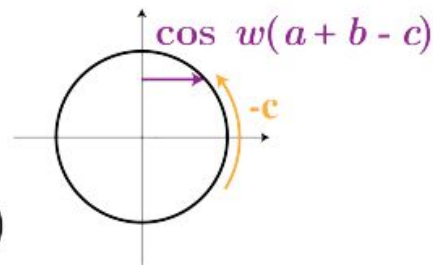
arXiv:2301.05217v3 [cs.LG] 19 Oct 2023

progress measures ~ mechanistic interpretability

reverseengineering learned behaviors into their individual components

Computes logits using further trig identities:

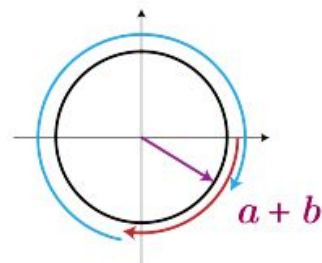
$$\begin{aligned}\text{Logit}(c) &\propto \cos(w(a + b - c)) \\ &= \cos(w(a + b)) \cos(wc) + \sin(w(a + b)) \sin(wc)\end{aligned}$$



Calculates sine and cosine of $a + b$ using trig identities:

$$\sin(w(a + b)) = \sin(wa) \cos(wb) + \cos(wa) \sin(wb)$$

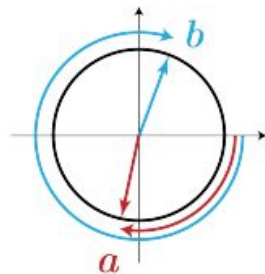
$$\cos(w(a + b)) = \cos(wa) \cos(wb) - \sin(wa) \sin(wb)$$



Translates one-hot a , b to Fourier basis:

$$a \rightarrow \sin(wa), \cos(wa)$$

$$b \rightarrow \sin(wb), \cos(wb)$$



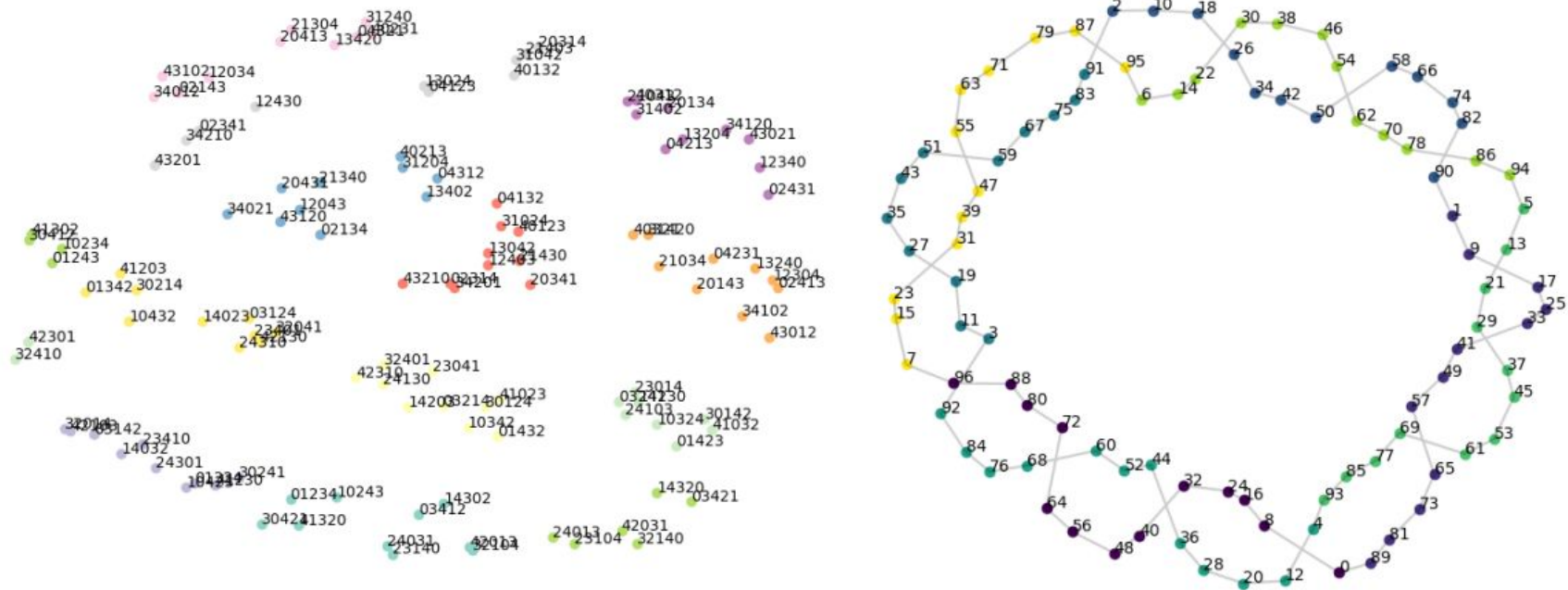
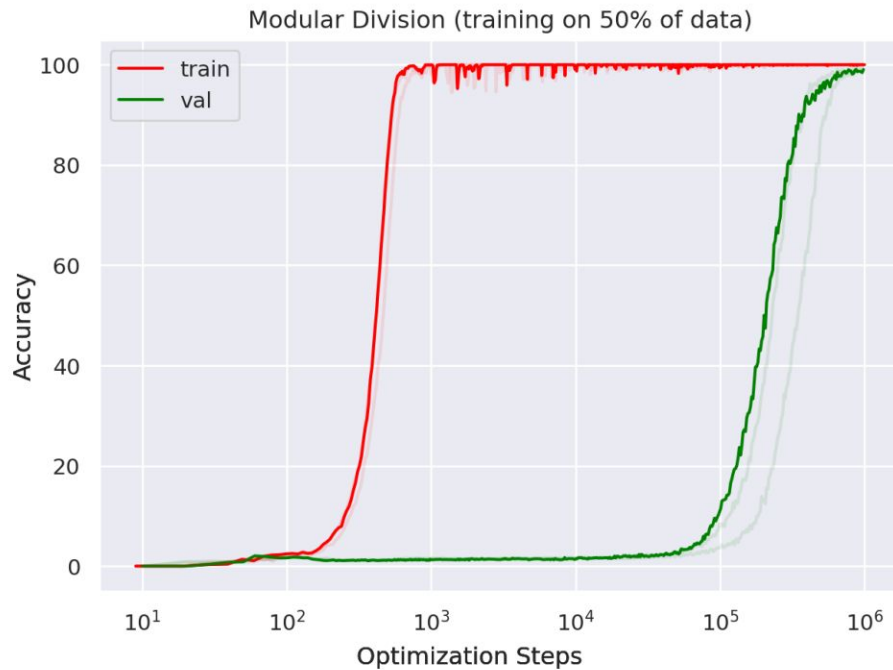


Figure 3: **Left.** t-SNE projection of the output layer weights from a network trained on S_5 . We see clusters of permutations, and each cluster is a coset of the subgroup $\langle (0, 3)(1, 4), (1, 2)(3, 4) \rangle$ or one of its conjugates. **Right.** t-SNE projection of the output layer weights from a network trained on modular addition. The lines show the result of adding 8 to each element. The colors show the residue of each element modulo 8.

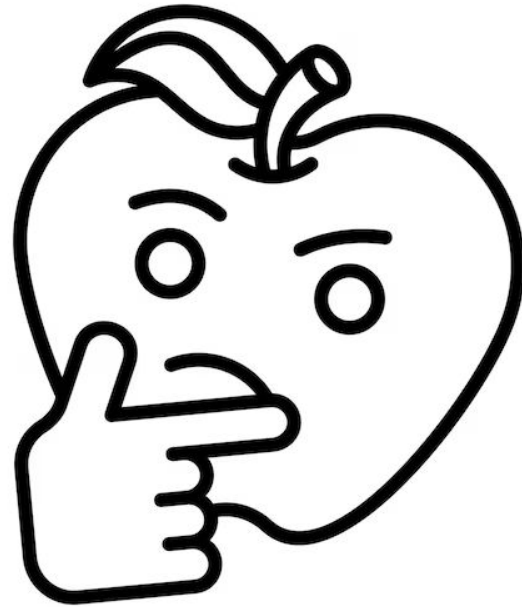
- **memorization**
- **circuit formation**
- **cleanup**

Cleanup = Grokking



Grokking = delayed generalization

Too Artificial?

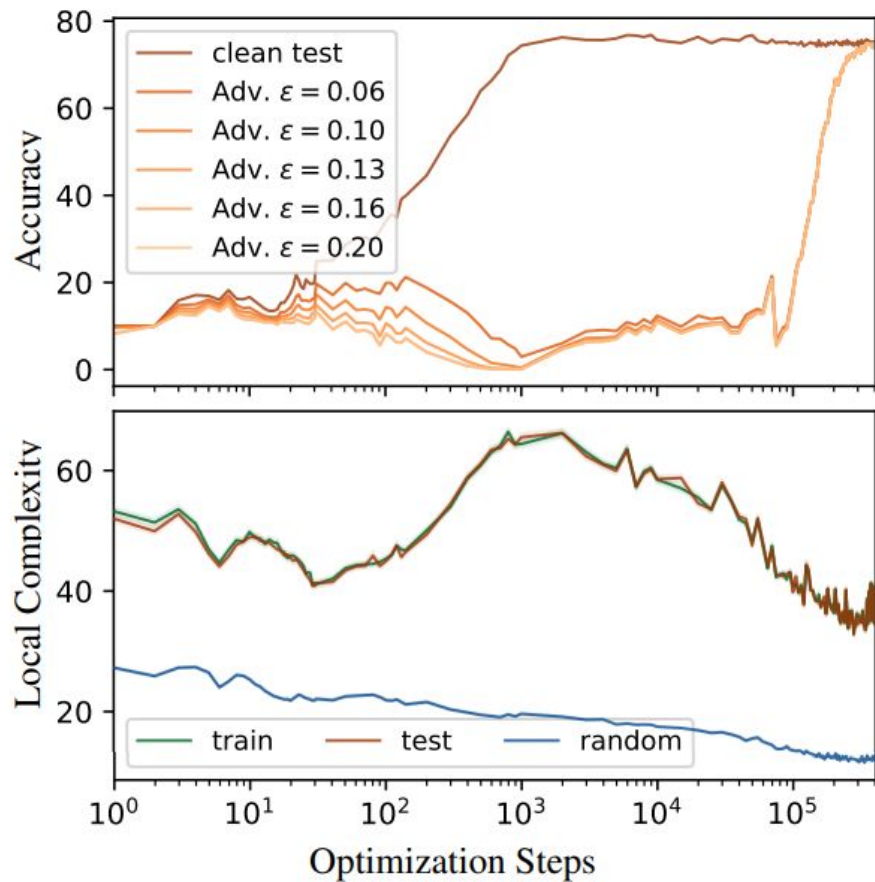


Deep Networks Always Grok and Here is Why

Ahmed Imtiaz Humayun¹ Randall Balestriero² Richard Baraniuk¹

arXiv:2402.15555v2 [cs.LG] 6 Jun 2024

- CNN on CIFAR10
- Resnet on Imagenette



Conclusion

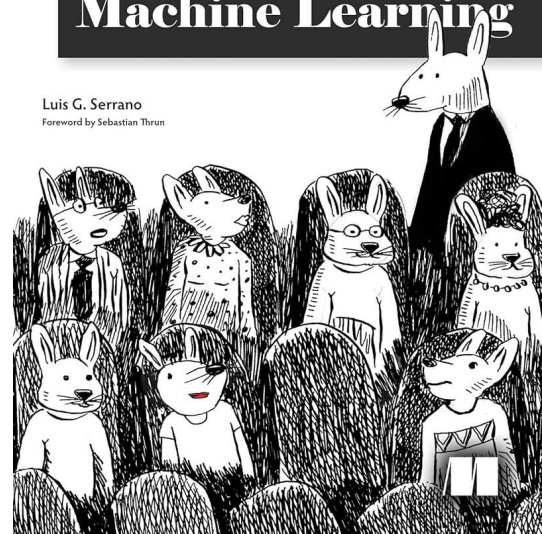
- **small datasets - ?**
- **grokking: accuracy ↗**
- **grokking: delayed generalization**
- **future - ?**



Thank you for attention!

Machine Learning

Luis G. Serrano
Foreword by Sebastian Thrun





Any questions?

