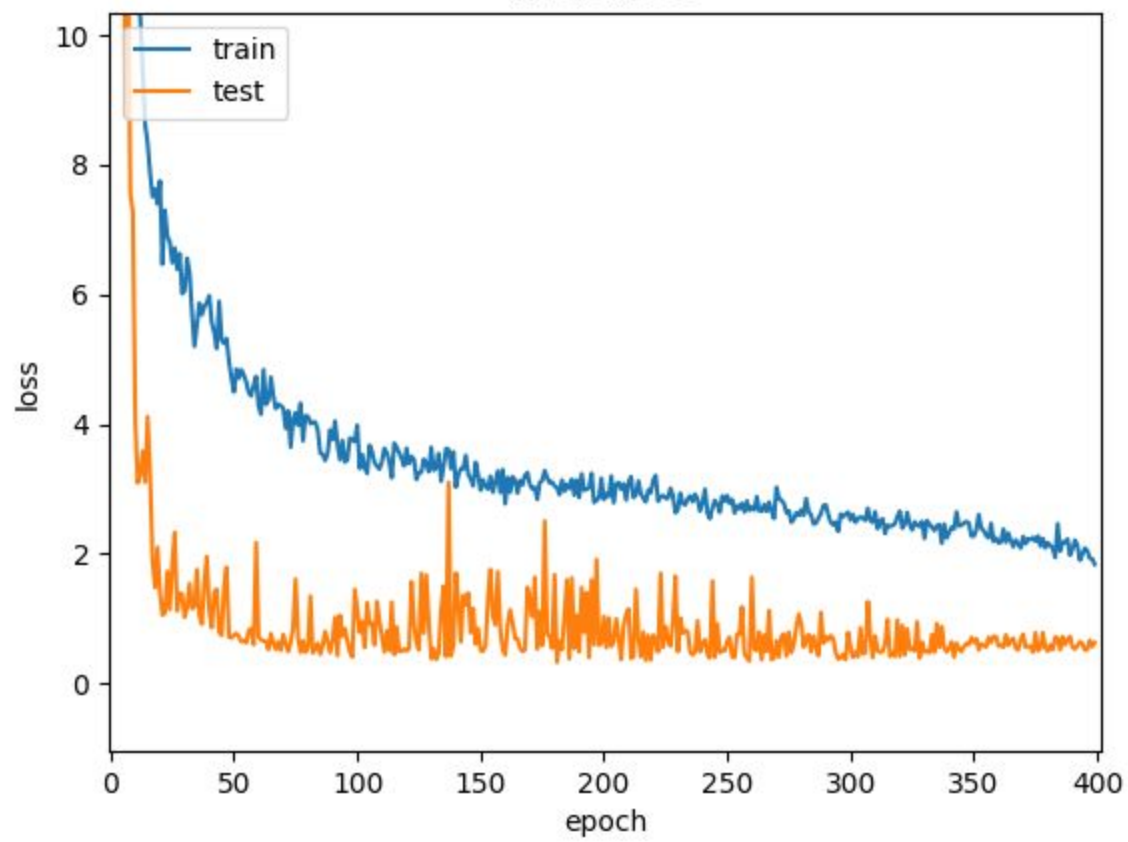


Grokking

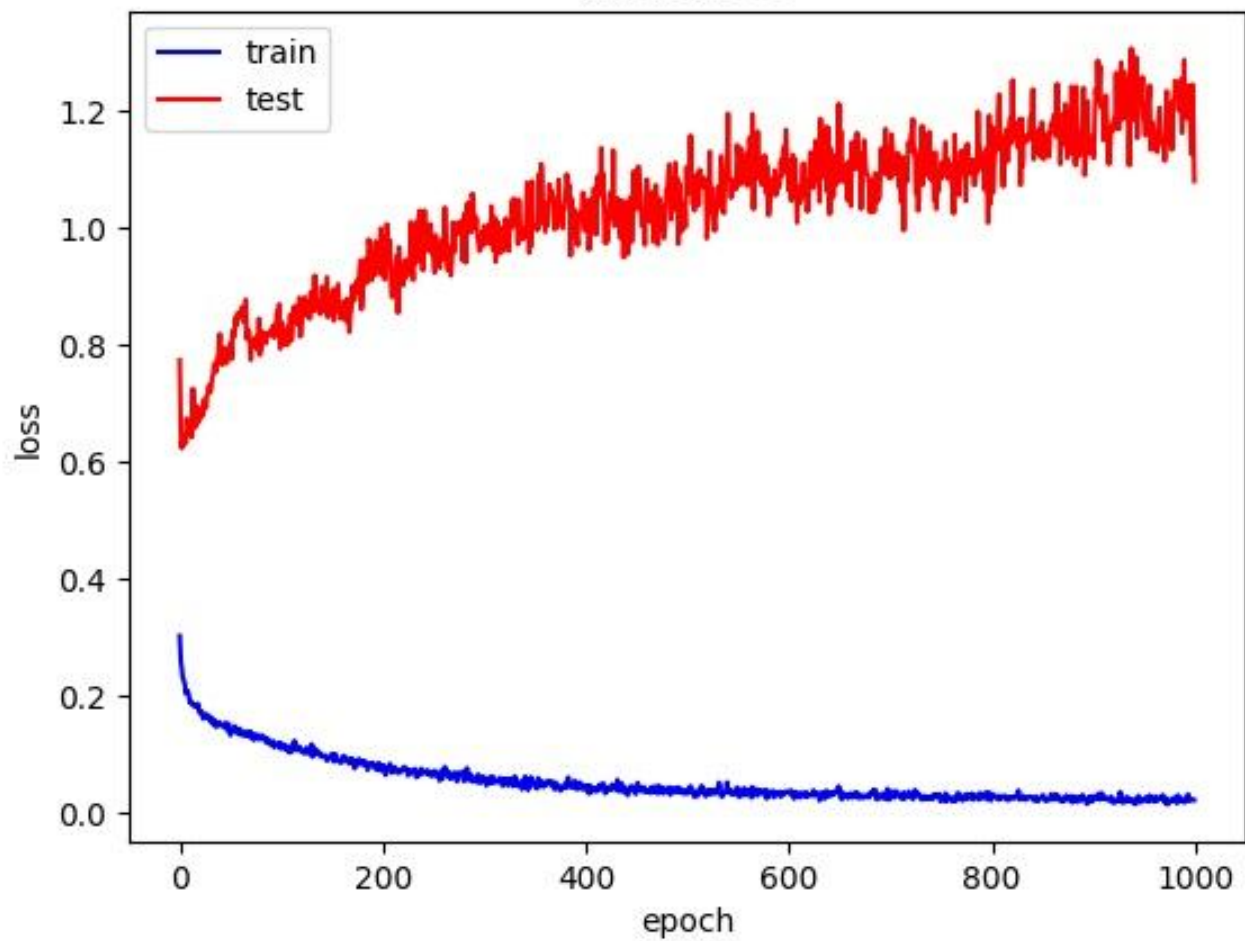
Михаил Доманин

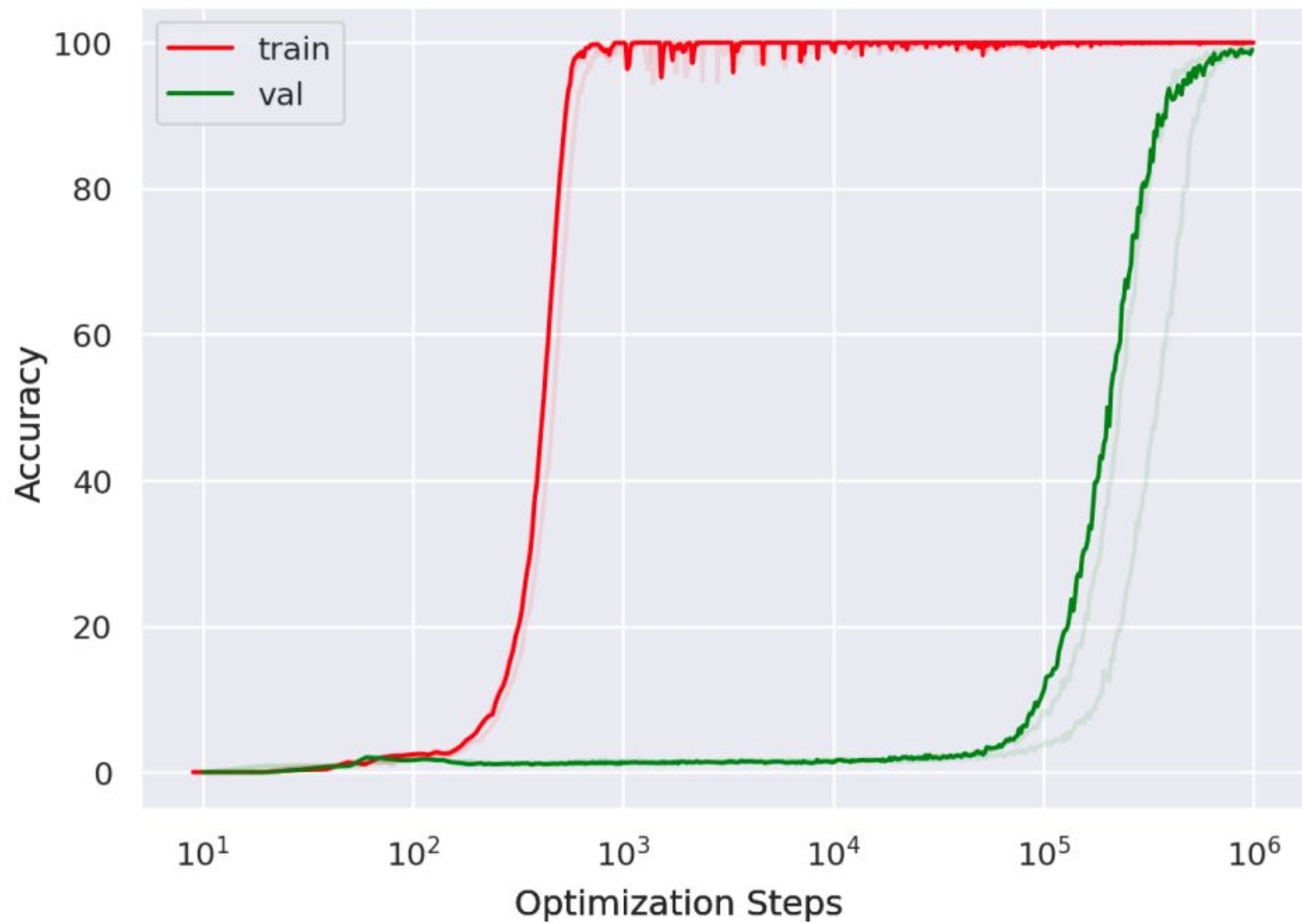
Grok – понимать что-либо полностью, буквально на интуитивном уровне, осознавая всю логику и закономерности, а не просто на уровне знания правил

model loss



model loss





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

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

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

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

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

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

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

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

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

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

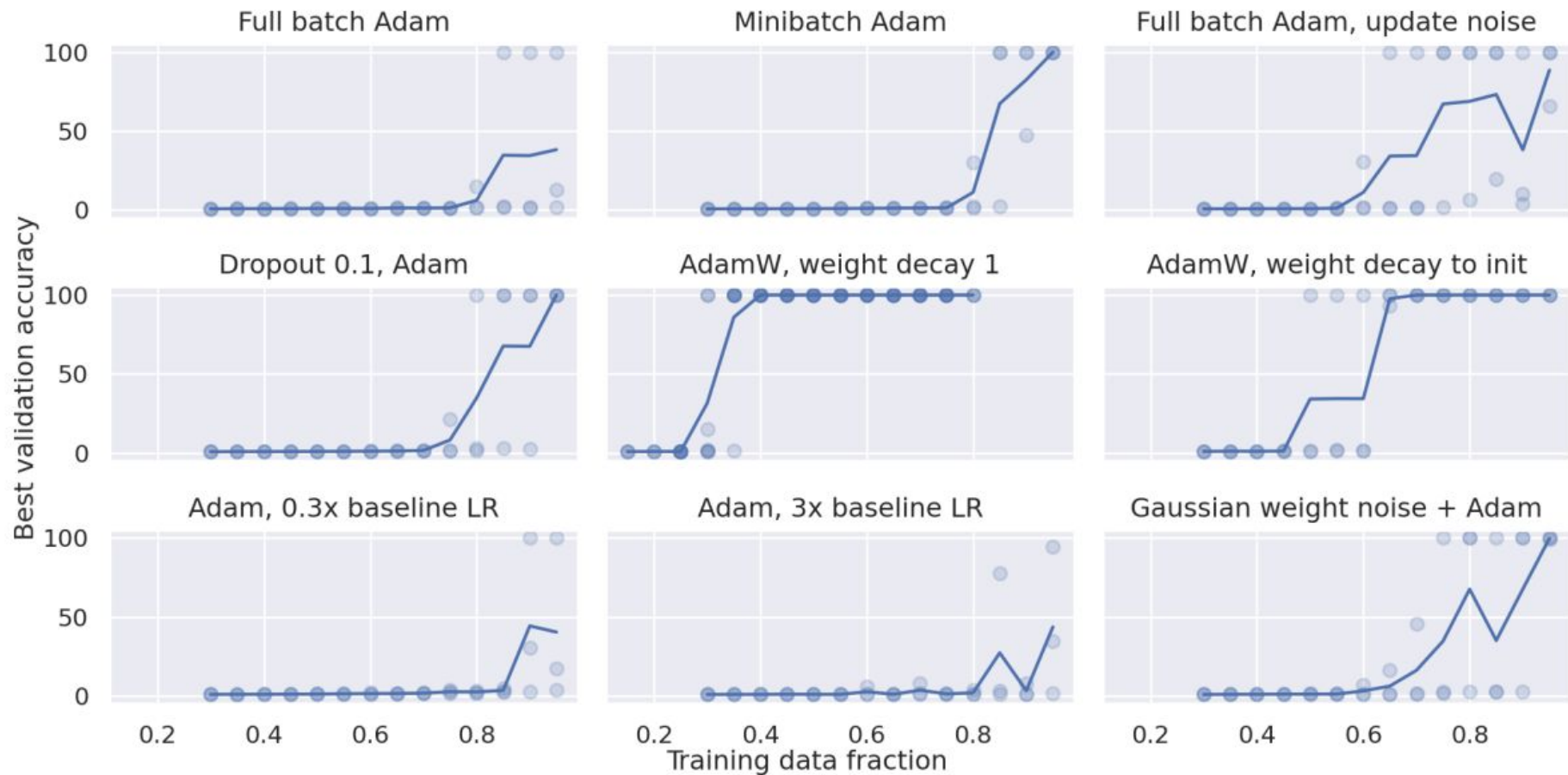
$$x \circ y = x \cdot y \text{ for } x, y \in S_5$$

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

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

Steps until generalization for product in abstract group S_5





The key feature of the algorithm is calculating $\cos(w(x + y)), \sin(w(x + y))$ with $w = \frac{2\pi}{p}k$ - this is a function of $x + y$ and be mapped to $x + y$, and because $\cos(wx)$ has period $\frac{p}{k}$ we get the $(\bmod p)$ part for free.

More concretely:

- Inputs x, y are given as one-hot encoded vectors in \mathbb{R}^p
- Calculates $\cos(wx), \cos(wy), \sin(wx), \sin(wy)$ via a Discrete Fourier Transform (This sounds complex but is just a change of basis on the inputs, and so is just a linear map)
 - $w = \frac{2\pi}{p}k$, k is arbitrary, we just need period dividing p
- Calculates $\cos(wx) \cos(wy), \cos(wx) \sin(wy), \sin(wx) \cos(wy), \sin(wx) \sin(wy)$ by multiplying pairs of waves in x and in y
- Calculates $\cos(w(x + y)) = \cos(wx) \cos(wy) - \sin(wx) \sin(wy)$ and $\sin(w(x + y)) = \sin(wx) \cos(wy) + \cos(wx) \sin(wy)$ by rearranging and taking differences
- Calculates $\cos(w(x + y - z)) = \cos(w(x + y)) \cos(wz) + \sin(w(x + y)) \sin(wz)$ via a linear map to the output logits z
 - This has an argmax at $z \equiv x + y \pmod{p}$, so post softmax we're done!

Гипотезы

- Grokking возникает в результате длительного блуждания модели по оптимумам функции потерь

Гипотезы

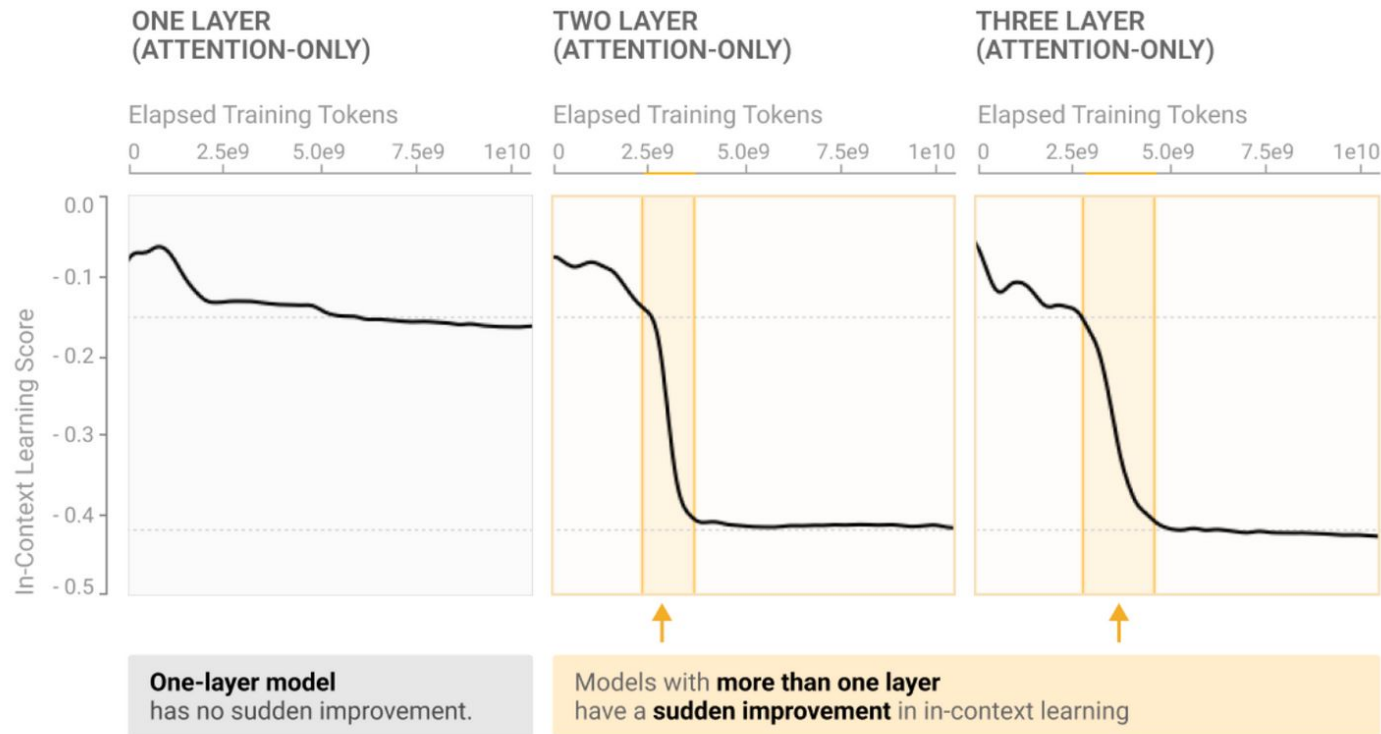
- Grokking возникает в результате длительного блуждания модели по оптимумам функции потерь
- Выход нейросети – комбинация значений разных схем. Некоторые из них потенциально полезны для уменьшения лоса, другие нет. Градиентный спуск “погасит” явно не важные признаки (обнуляя коэффициенты перед ними) и усилит полезные, то есть модель будет постепенно все ближе и ближе к тому, чтобы понять истинную закономерность.



Входные данные

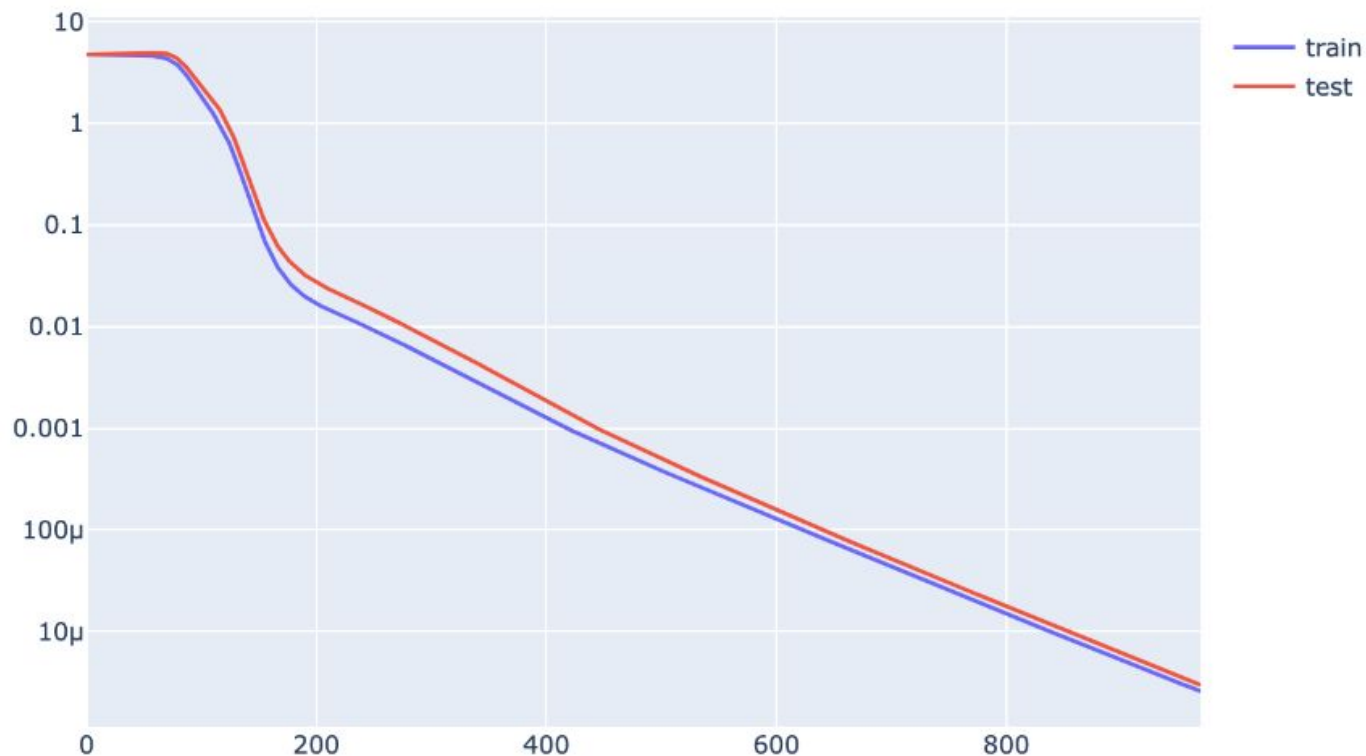
- Данных слишком мало — модель легко запоминает, нет стимула к обобщению
- Данных слишком много — модель легко обобщает
- Данных ни много, ни мало — проявляется Grokking

MODELS WITH MORE THAN ONE LAYER HAVE AN ABRUPT IMPROVEMENT IN IN-CONTEXT LEARNING



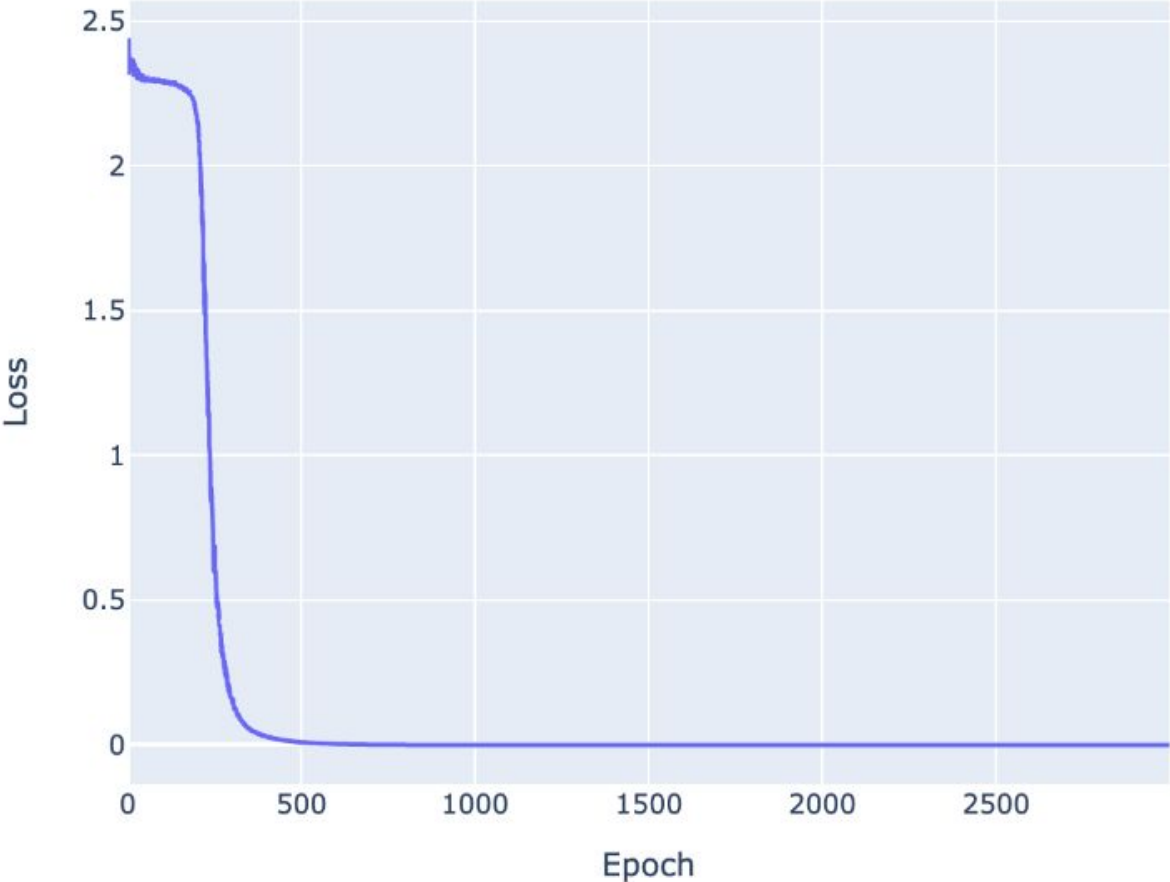
— We highlight the “**phase change**” period of training in plots to make visual comparison between plots easier. The highlighted region is selected for each model based on the derivative of in-context learning.

Train + Test Loss curves for modular addition trained on 95% of the data

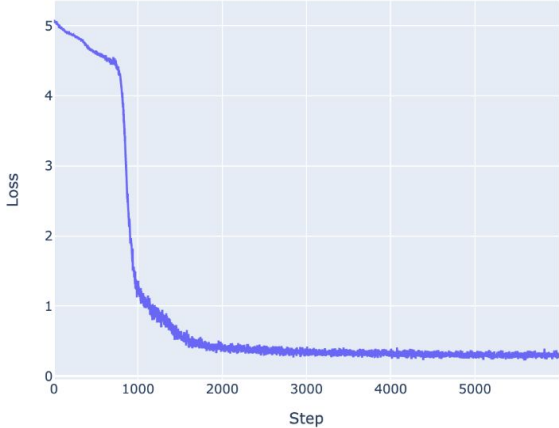


Modular addition mod 113 loss curve, trained on 95% of the data

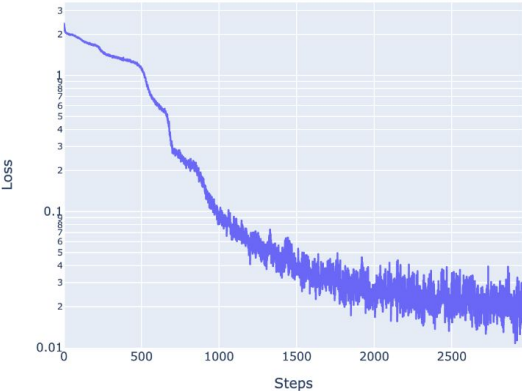
Skip Trigram Loss Infinite Data (Linear Scale)



Repeated Subsequence Prediction Infinite Data Training



Phase Change in 5 Digit Addition Infinite Data Training Curve



Grokking = Регуляризация + Фазовые изменения + Ограниченные данные

Что насчет реальных задач?

- Grokking полагается на интерпретируемость модели
- Шумы сильно мешают возникновению эффекта
- Сложные модели – комбинация большого количества схем, поиск обобщающего решения вычислительно может быть очень долгим
- Однако...

Per-digit Loss Curves for 5 digit addition (Infinite Data)

