# Large Language Models

## Scaling Laws

ELL881 · AIL821

**Sourish Dasgupta**
Assistant Professor, DA-IICT, Gandhinagar
*https://daiict.ac.in/faculty/sourish-dasgupta*

# Kaplan Scaling Laws at a glance:

| Power Law | Scale (tokenization-dependent) |
|---|---|
| $\alpha_N = 0.076$ | $N_c = 8.8 \times 10^{13}$ params (non-embed) |
| $\alpha_D = 0.095$ | $D_c = 5.4 \times 10^{13}$ tokens |
| $\alpha_C = 0.057$ | $C_c = 1.6 \times 10^7$ PF-days |
| $\alpha_C^{min} = 0.050$ | $C_c^{min} = 3.1 \times 10^8$ PF-days |
| $\alpha_B = 0.21$ | $B_* = 2.1 \times 10^8$ tokens |
| $\alpha_S = 0.76$ | $S_c = 2.1 \times 10^3$ steps |

| Parameters | Data | Compute | Batch Size | Equation |
|---|---|---|---|---|
| $N$ | $\infty$ | $\infty$ | Fixed | $L(N) = (N_c/N)^{\alpha_N}$ |
| $\infty$ | $D$ | Early Stop | Fixed | $L(D) = (D_c/D)^{\alpha_D}$ |
| Optimal | $\infty$ | $C$ | Fixed | $L(C) = (C_c/C)^{\alpha_C}$ (naive) |
| $N_{opt}$ | $D_{opt}$ | $C_{min}$ | $B \ll B_{crit}$ | $L(C_{min}) = (C_c^{min}/C_{min})^{\alpha_C^{min}}$ |
| $N$ | $D$ | Early Stop | Fixed | $L(N,D) = \left[ \left(\frac{N_c}{N}\right)^{\frac{\alpha_N}{\alpha_D}} + \frac{D_c}{D} \right]^{\alpha_D}$ |
| $N$ | $\infty$ | $S$ steps | $B$ | $L(N,S) = \left(\frac{N_c}{N}\right)^{\alpha_N} + \left(\frac{S_c}{S_{min}(S,B)}\right)^{\alpha_S}$ |

|| Joint.

# Is there any other alternative law?

# Turns out there is!

*Lower bound*

$$\text{Loss}(N_T, D) = \frac{N_c}{N_T^\alpha} + \frac{D_c}{D^\beta} + E,$$

$$\text{Loss}(N_T, C_T) = \frac{N_c}{N_T^\alpha} + \frac{D_c}{(C_T/6N_T)^\beta} + E$$

*Total parameters.*

*Loss that is intrinsic to the Dataset.*

**Chinchilla (Hoffman) Scaling Law**

# The Chinchilla (Hoffman) Scaling Law

$$\text{Loss}(N_T, D) = \frac{N_c}{N_T^\alpha} + \frac{D_c}{D^\beta} + E \longrightarrow L(N, D) = 1.69 + \frac{406.4}{N^{0.34}} + \frac{410.7}{D^{0.28}}$$

$$N_{opt}(C) = G(C/6)^a \qquad D_{opt}(C) = G^{-1}(C/6)^b$$

$$\text{where} \quad G = \left(\frac{\alpha A}{\beta B}\right)^{\frac{1}{\alpha+\beta}} \qquad a = \frac{\beta}{\alpha + \beta} \qquad b = \frac{\alpha}{\alpha + \beta}$$

Fitting the constants, yields: $\alpha \approx \beta$ $\rightarrow$ *best curve fitting*

i.e. equal scaling of **N** and **D**.

# Chinchilla Scaling Law vs. Kaplan Scaling Law

$$\text{Kaplan:} N^*_{\backslash E} \propto C^{0.73}_{\backslash E}$$

$$\text{Chinchilla:} N^*_T \propto C^{0.50}_T.$$

*arbour
Embed Matrix.*

Performance penalty is $N^{0.75} / D$
- if model increases 8x, dataset must increase 5x

VS.

Fitting the constants, yields: $\alpha \approx \beta$

i.e. equal scaling of **N** and **D**.

*param for Embed Matrix.*

$$N_T = N_E + N_{\backslash E}, \qquad C_T = 6N_T D = 6(N_E + N_{\backslash E})D,$$

$$N_E = (h + v)d, \qquad C_{\backslash E} = 6N_{\backslash E}D.$$

*context . vocabs -*

# The (revised) Chinchilla Scaling Law

$$L(N, D) = 1.82 + \frac{514.0}{N^{0.35}} + \frac{2115.2}{D^{0.37}}$$

*Scaling factor by Dahoost.*

→ better fit.

$$L(N, D) = 1.69 + \frac{406.4}{N^{0.34}} + \frac{410.7}{D^{0.28}}$$

→ Chinchilla law.

# Is it a problem with our point-of-*view*?

# LLMs "*seems*" to get more intelligent with the following:



Amount of training data
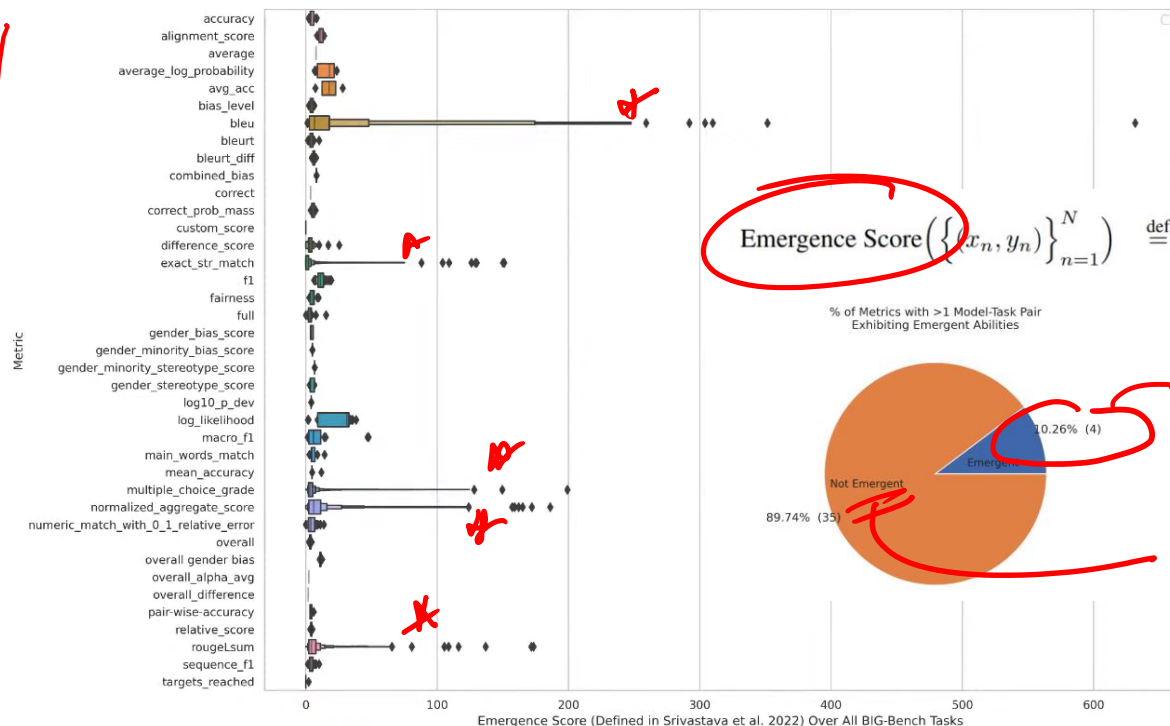
Model size (parameters)

Amount of compute (or time)
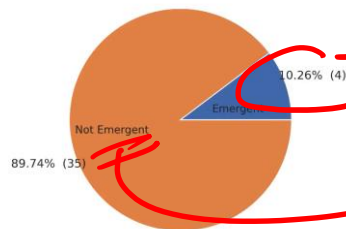
$$\mathcal{L}_{CE}(N) = \left(\frac{N}{c}\right)^{\alpha}$$

# Motivation: Not all metrics score same (Emergence Score)

$$\text{Emergence Score}\left(\left\{(x_n, y_n)\right\}_{n=1}^{N}\right) \overset{\text{def}}{=} \frac{\text{sign}(\arg\max_i y_i - \arg\min_i y_i)(\max_i y_i - \min_i y_i)}{\sqrt{\text{Median}(\{(y_i - y_{i-1})^2\}_i)}}$$

% of Metrics with >1 Model-Task Pair
Exhibiting Emergent Abilities

10.26% (4)

Emergent

89.74% (35)

Not Emergent

Emergence Score (Defined in Srivastava et al. 2022) Over All BIG-Bench Tasks

# Is your accuracy metric non-linear or discontinuous?

*benchmark.*

> 92% of BIG-BENCH:

*matters*

$$\text{Multiple Choice Grade} \overset{\text{def}}{=} \begin{cases} 1 & \text{if highest probability mass on correct option} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Exact String Match} \overset{\text{def}}{=} \begin{cases} 1 & \text{if output string exactly matches target string} \\ 0 & \text{otherwise} \end{cases}$$
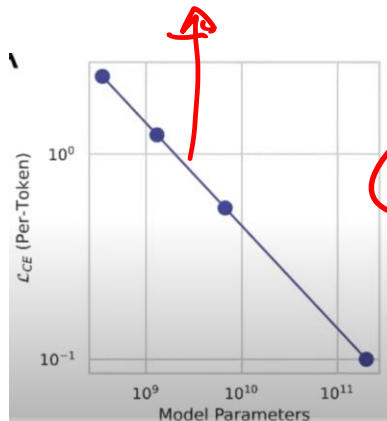


16.42% (11)
exact_str_match
bleu 2.99% (2)
normalized_aggregate_score
4.48% (3)
multiple_choice_grade
76.12% (51)

*Too challenging for smaller models!*
*Is it really worth??*

# Power Law in play!



$$\mathcal{L}_{CE}(N) = \left(\frac{N}{c}\right)^{\alpha} = -\log \hat{p}_{v^*}(N)$$

*Kaplan*

$\mathcal{L}()$ CROSS ENTROPY

$$\hat{p}_{v^*}(N) = \exp\left(-(N/c)^{\alpha}\right)$$

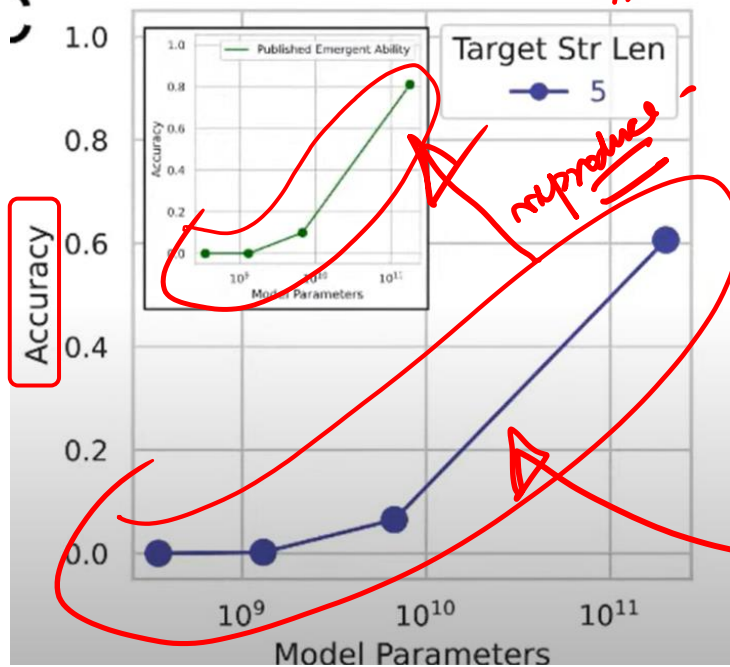$$L(D_0, N) = B_0 + \frac{B_1}{N^{\alpha_w}} + \dots$$

Length of token

S-curve.

# Problem with Non-linear Measure: Eg.: *Exact string match*

**Task**: *Add k-digit integers*



| 1 | if all K+1 digits in model's output are correct |
| 0 | otherwise |

$$\hat{p}_{v^*}(N) = \exp\left(-(N/c)^{\alpha}\right)$$

$$\text{Accuracy}(N) \approx p_N(\text{single token correct})^{\text{num. of tokens}}$$

# Change of perspective: Measure: *Edit distance*

**Task**: *Add k-digit integers*



| 1 | if all K+1 digits in model's output are correct |
|---|---|
| 0 | otherwise |

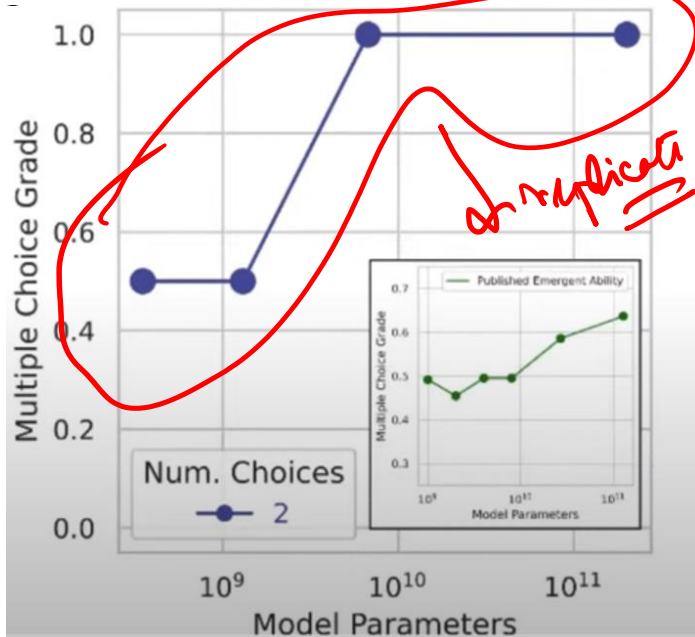$$\hat{p}_{v^*}(N) = \exp\left(-(N/c)^\alpha\right)$$

$$\text{Edit Distance}(N) \approx L\left(1 - p_N(\text{single token correct})\right)$$

# Problem with ==Discontinuous== Measure: Eg.: *MCG*

**Task**: *Choose one of two*



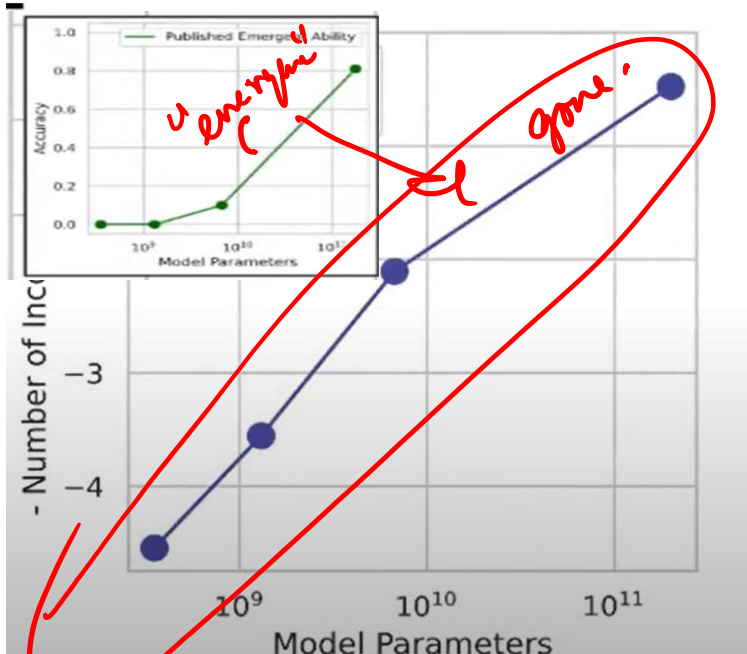| 1 | if highest probability mass on correct option |
|---|---|
| 0 | otherwise |

$$\hat{p}_{v^*}(N) = \exp\left(-(N/c)^{\alpha}\right)$$

# Change of perspective: Measure: *Brier Score*

**Task**: *Choose one of two*



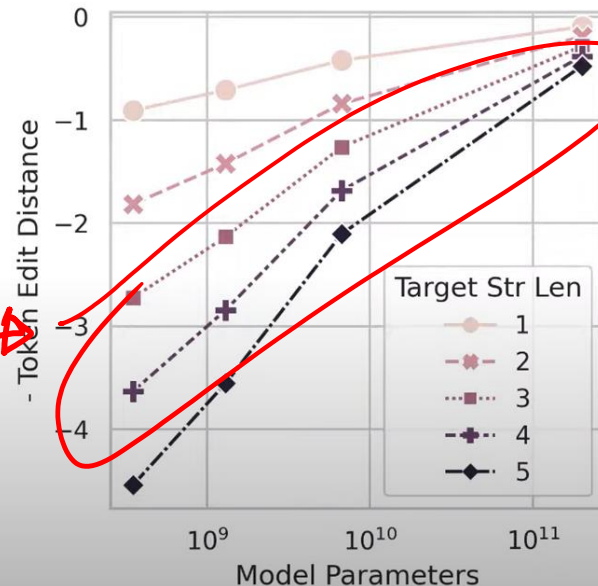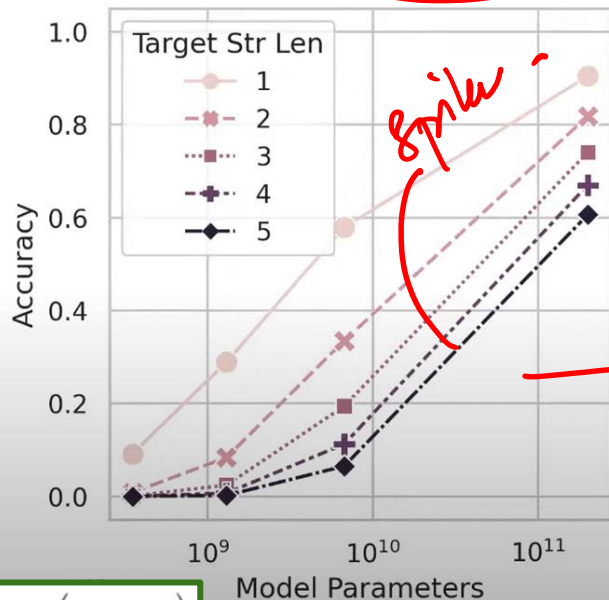| 1 | if all K+1 digits in model's output are correct |
|---|---|
| 0 | otherwise |

$$\hat{p}_{v^*}(N) = \exp\left(-(N/c)^\alpha\right)$$

Brier Score = (1 – probability mass on correct option)$^2$

# Prediction: Power Law vs. Near-Linear counterpart
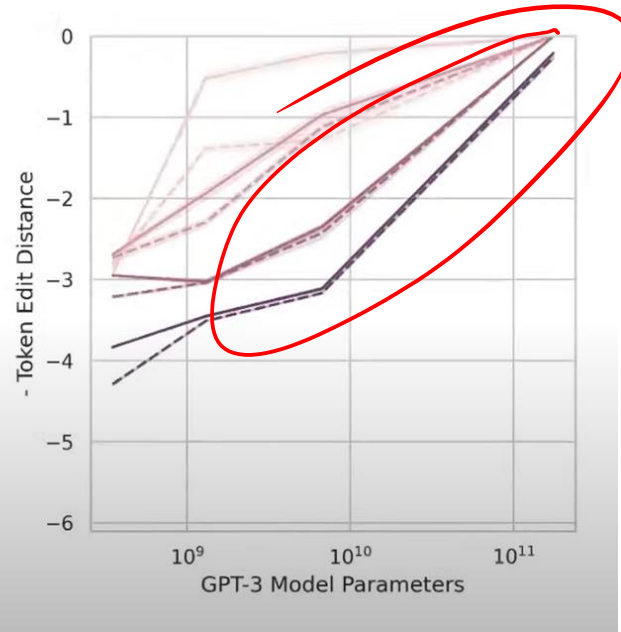


$$\hat{p}_{v^*}(N) = \exp\left(-(N/c)^\alpha\right)$$

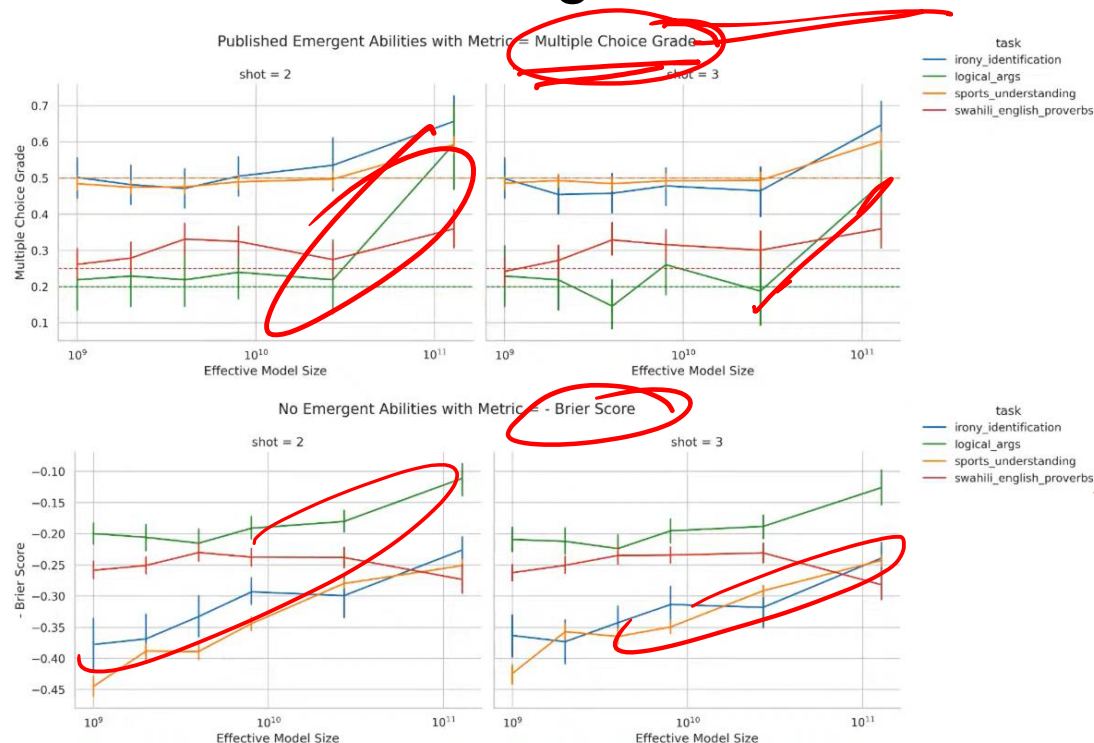# Results on GPT3.5/3: Task: 2-digit integer multiplication



$$\hat{p}_{v^*}(N) = \exp\left(-(N/c)^\alpha\right)$$

# Does the claim work for Google BIG-BENCH benchmark?

# Key Takeaways

- Want to predict <mark>without the theatrics</mark>? Choose a _metric that's "soft"_ (in the continuous sense)

- There's _no sudden jump_ in reality ("_most_" can be predicted on a near-linear scale)

- Do  we really need the power law of scale? Maybe not!