

317 **A Approximate Behavior of Metrics on Sequential Data**

318 How do different metrics behave when used to measure autoregressive model outputs? Precisely
 319 answering this question is tricky and possibly analytically unsolvable, so we provide an approximate
 320 answer here.

321 Notationally, we consider N test data of length L (here, length is measured in tokens) with tar-
 322 gets denoted $t_n \stackrel{\text{def}}{=} (t_{n1}, t_{n2}, \dots, t_{nL})$, the autoregressive model has a true-but-unknown per-token er-
 323 ror probability of $\epsilon \in [0, 1]$ and the model outputs prediction $\hat{t}_n \stackrel{\text{def}}{=} (\hat{t}_{n1}, \hat{t}_{n2}, \dots, \hat{t}_{nL})$. This assumes
 324 that the model’s per-token error probability is constant, which is empirically false, but modeling the
 325 complex dependencies of errors is beyond our scope.

326 **A.1 Per-Token Error Probability is Resolution-Limited**

327 Note that because we have N test data, each of length L , our resolution for viewing the per-token
 328 error probability ϵ is limited by $1/NL$. Here, resolution refers to “the smallest interval measurable
 329 by a scientific instrument; the resolving power.” To explain what resolution means via an example,
 330 suppose one wants to measure a coin’s probability of yielding heads. After a single coin flip, only
 331 two outcomes are possible (H, T), so the resolution-limited probability of heads is either 0 or 1. After
 332 two coin flips, four outcomes are possible (HH, HT, TH, TT), so the resolution-limited probability
 333 of heads is now one of 0, 0.5, 1. After F coin flips, we can only resolve the coin’s probability of
 334 yielding heads up to $1/F$. Consequently, we introduce a resolution-limited notation:

$$\lfloor a \rfloor_b \stackrel{\text{def}}{=} a \text{ rounded to the nearest integer multiple of } 1/b \quad (3)$$

335 **A.2 Token Edit Distance**

336 We first consider an adaptation of the Levenshtein (string edit) distance for models that function
 337 on tokens rather than characters, an adaptation we term the *token edit distance*. The token edit
 338 distance between two token sequences t_n, \hat{t}_n is defined as the integer number of additions, deletions
 339 or substitutions necessary to transform t_n into \hat{t}_n (or vice versa).

$$\text{Token Edit Distance}(t_n, \hat{t}_n) \stackrel{\text{def}}{=} \text{Num Substitutions} + \text{Num. Additions} + \text{Num. Deletions} \quad (4)$$

$$= \sum_{\ell=1}^L \mathbb{I}[t_{n\ell} \neq \hat{t}_{n\ell}] + \text{Num. Additions} + \text{Num. Deletions} \quad (5)$$

$$\geq \sum_{\ell=1}^L \mathbb{I}[t_{n\ell} \neq \hat{t}_{n\ell}] \quad (6)$$

340 The expected token edit distance is therefore:

$$\mathbb{E}[\text{Token Edit Distance}(t_n, \hat{t}_n)] \geq \mathbb{E}\left[\sum_{\ell=1}^L \mathbb{I}[t_{n\ell} \neq \hat{t}_{n\ell}]\right] \quad (7)$$

$$= \sum_{\ell=1}^L p(t_{n\ell} \neq \hat{t}_{n\ell}) \quad (8)$$

$$\approx L(1 - \epsilon) \quad (9)$$

341 The resolution-limited expected token edit distance is therefore:

$$\lfloor \mathbb{E}[\text{Token Edit Distance}(t_n, \hat{t}_n)] \rfloor_{NL} \geq L \left(1 - \lfloor \epsilon \rfloor_{NL}\right) \quad (10)$$

342 From this, we see that the expected token edit distance scales approximately linearly with the
 343 resolution-limited per-token probability. The real rate is slightly higher than linear because addi-
 344 tions and deletions contribute an additional non-negative cost, but modeling this requires a model
 345 of how likely the model is to overproduce or underproduce tokens, which is something we do not
 346 currently possess.

347 A.3 Accuracy

$$\text{Accuracy}(t_n, \hat{t}_n) \stackrel{\text{def}}{=} \mathbb{I}[\text{No additions}] \mathbb{I}[\text{No deletions}] \prod_{l=1}^L \mathbb{I}[t_{nl} = \hat{t}_{nl}] \quad (11)$$

$$\approx \prod_{l=1}^L \mathbb{I}[t_{nl} = \hat{t}_{nl}] \quad (12)$$

348 As with the Token Edit Distance (App. A.3), we ignore how likely the language model is to over-
 349 produce or underproduce tokens because we do not have a good model of this process. Continuing
 350 along,

$$\mathbb{E}[\log \text{Accuracy}] = \sum_l \mathbb{E}[\log \mathbb{I}[t_{nl} = \hat{t}_{nl}]] \quad (13)$$

$$\leq \sum_l \log \mathbb{E}[\mathbb{I}[t_{nl} = \hat{t}_{nl}]] \quad (14)$$

$$\approx L \log(1 - \epsilon) \quad (15)$$

351 Taking an approximation that would make most mathematicians cry:

$$\mathbb{E}[\text{Accuracy}] \approx \exp(\mathbb{E}[\log \text{Accuracy}]) \quad (16)$$

$$= (1 - \epsilon)^L \quad (17)$$

$$(18)$$

352 This reveals that accuracy **approximately** falls geometrically with target token length. The
 353 resolution-limited expected accuracy is therefore:

$$\lfloor \mathbb{E}[\text{Accuracy}] \rfloor_{NL} = \lfloor (1 - \epsilon)^L \rfloor_{NL} \quad (19)$$

354 From this we can see that choosing a nonlinear metric like Accuracy is affected significantly more
 355 by limited resolution because Accuracy forces one to distinguish quantities that decay rapidly.

356 A.4 ROUGE-L-Sum

357 Another BIG-Bench metric [28] is ROUGE-L-Sum [23], a metric based on the longest common
 358 subsequence (LCS) between two sequences. Section 3.2 of [23] gives the exact definition, but the
 359 key property is that ROUGE-L-Sum measures the “union” LCS, which means “stitching” together
 360 LCSs across the candidate and multiple references. As explained in the original paper: if the candi-
 361 date sequence is $c = w_1 w_2 w_3 w_4 w_5$, and if there are two reference sequences $r_1 = w_1 w_2 w_6 w_7 w_8$
 362 and $r_2 = w_1 w_3 w_8 w_9 w_5$, then $LCS(r_1, c) = w_1 w_2$ and $LCS(r_2, c) = w_1 w_3 w_5$, then the *union*
 363 -LCS of c, r_1, r_2 is $w_1 w_2 w_3 w_5$, with length 4. Intuitively, this disproportionately benefits models
 364 with smaller error rates because their mistakes can be “stitched” across multiple references; this is
 365 confirmed in simulation (Fig. 9).

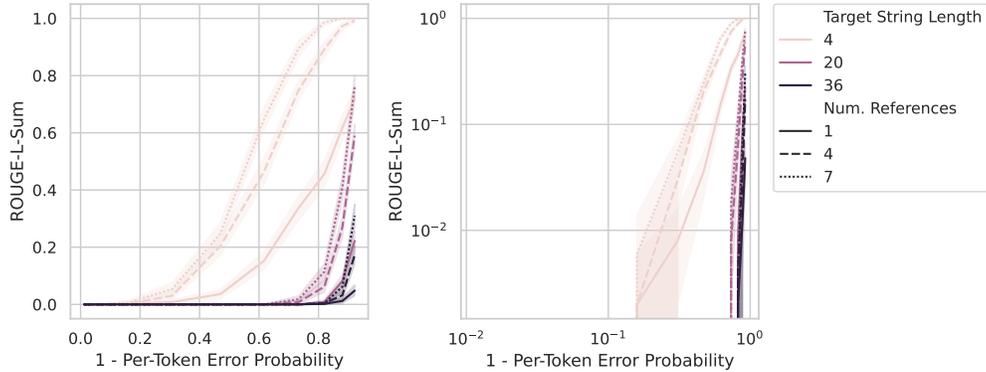


Figure 9: **ROUGE-L-Sum is a sharp metric.** Simulations show that as the per-token error probability slightly increase (e.g. from 0.05 to 0.1), the ROUGE-L-Sum metric sharply falls.

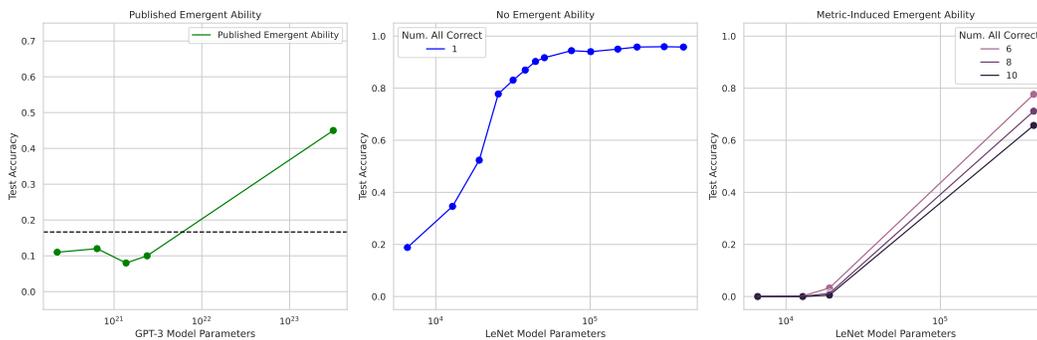


Figure 10: **Induced emergent MNIST classification ability in convolutional networks.** (A) A published emergent ability from the BIG-Bench Grounded Mappings task [33]. (B) LeNet trained on MNIST [21] displays a predictable, commonplace sigmoidal increase in test accuracy as model parameters increase. (C) When accuracy is redefined as correctly classifying K out of K independent test data, this newly defined metric induces a seemingly unpredictable change.

366 B Inducing Emergent Abilities in Networks on Vision Tasks

367 B.1 Emergent Classification of MNIST Handwritten Digits by Convolutional Networks

368 We begin by inducing an emergent classification ability in a LeNet convolutional neural network
 369 family [22], trained on the MNIST handwritten digits dataset [21]. This family displays smoothly
 370 increasing test accuracy as the number of parameters increase (Fig. 10B). To emulate the accuracy
 371 metric used by emergence papers [8, 33, 28], we use *subset accuracy*: 1 if the network classifies K
 372 out of K (independent) test data correctly, 0 otherwise. Under this definition of accuracy, the model
 373 family displays an “emergent” ability to correctly classify sets of MNIST digits as K increases from
 374 1 to 5, especially when combined with sparse sampling of model sizes (Fig. 10C). This convolutional
 375 family’s emergent classification ability qualitatively matches published emergent abilities, e.g., at
 376 the BIG-Bench Grounded Mappings task [33] (Fig. 10A).