# A Examples of Glitch Tokens

In this section, we present a subset of special tokens identified in the DeepSeek-V3-0324 model. Due to space constraints, we focus on glitch tokens composed solely of English alphabetic characters.

Note that the leading '‿' in tokens such as ‿CreatureTPL indicates a leading space.

**Repeat-Type Glitch Tokens:**

‿kinabugnawan; ugnawan; ‿kinainitan; ‿giiniton; ebabkan; eredReader; Beskjeftigelse; ‿ngalan; inheritdoc; ‿gihabogon; ‿CreatureTPL; jeftigelse; ‿nM; ahabogang; ‿zituzten; ‿mediabestanden; Koordenatuak; sweise; ordenatuak; ‿anhianhi; ‿zeuden; dfunding; Administrazioa; Tallennettuna; eltemperaturen; ennettuna; armaceut; ‿kasadpan; Siyentipikinhong; ‿kahaboga; Siyentipik; ‿nahimutangan; ‿nahimut; Kaliwatan; ‿Noruwega; ‿amihanan; ‿habagatan; asarangang; ‿rozpoc; ‿kahenera; asilkan; ‿kinaugahan; ‿kinabasaan; ‿kinaug; ‿burujabe; ‿kabanay; Kasipak; eredWriter; ‿kinahabogang; ‿Nameeee; ‿numbersaplenty; ‿mPa; ‿factorisate; ‿kasarangang; ‿Ginhadi; ‿everydaycalculation; nig; asterxml; ‿nalukop; ‿pagklasipika; ‿pagklas; hematica; adrado; ‿Tsiahy; ‿hilabihan; ‿Substantiivi; ‿MentionsView; HasColumnType; ‿kasarigan; ‿Kalkulado; tanler; Siyentipiko; ‿Nahimutang; ‿matoanteny; ‿ulohan; ‿udalerria; ‿matoant; ‿DelaLika; ‿DelaL; Oppslagsverk; itetsdata; ulagway; Azalera; ‿itandi; dasarkan; ‿kinadul; ‿Kokoteksti; ftigelse; Dentsitatea; ellationToken; alohany; ‿kontsultatua; ultatua; iyembre; engono; asadpan; ‿basihan; ‿kabungtoran; ‿kabungtor; asadpang; ‿einges; idisciplinary; Nig; ‿Tinipong; owanych; ‿frantsay; ahimut; ‿Gikuha; izacji; ‿sidlakan; ‿talagsaon; ‿talags; ‿jednine; ‿suicide; ebizitza; Tiganos; superscriptsubscript; Kadaghanon; Kadaghan; ‿Pagbuok; ‿Britonhon; ‿Nieukerken; ‿kainiton; Kasarangang; kowych; Siyent; ‿nahilalakip; tanleria; nisone; ‿munisipyo; ‿iombonana; ‿Erreferentziak; visiae; ‿habagatang; ytocin; referentziak; ricula; delete; entsitatea; ‿sabwag; ‿kondado; Kahanay; Kahutong; Kaginharian; ‿Pagklasipikar; ‿Pagklas; Kabanay; onimy; Kahenera; ‿Viitattu; inharian; ret; ‿pakigbingkil; Kapunoang; bingkil; ‿Vertaisarvioitu; ‿Frantsay; aisarvioitu; Biztanleria; rapist; unoang; ‿ploraly; ahimutang; jonalitet; jektiv; ‿CategoryTreeLabel; itattu;

**Spell-Type Glitch Tokens:**

Longrightarrow; radation; rinsic; ‿kinabugnawan; ugnawan; ‿kinainitan; igtausend; itatea; ‿giiniton; ulinum; ebabkan; eredReader; ‿udaler; opsida; ‿niini; Beskjeftigelse; UIKit; ercase; ‿CreatureTPL; ‿gihabogon; jeftigelse; ‿palibot; ordable; orylation; iedenis; Bungtod; cohol; ruptedException; opausal; rebbero; ‿mediabestanden; ‿zituzten; bernetes; ahabogang; ‿nalista; Koordenatuak; niejsze; ‿anhianhi; chnung; sweise; ordenatuak; ‿zeuden; ‿ngOnInit; atieve; dfunding; ‿mosunod; Administrazioa; etCode; usztus; eltemperaturen; utnya; ‿daerah; yalgia; Hakutulos; Tallennettuna; ‿estekak; ennettuna; ispiele; yskland; rattutto; pskohrvatski; erequisite; ‿anglisy; otechnol; utterstock; ‿kasadpan; ziako; ‿ibabaw; rophot; ‿FullEDMFunc; EDMFunc; ‿adtong; ‿Noruwega; ‿nahimutangan; ‿nahimut; ‿kahaboga; abogon; Kaliwatan; Siyentipikinhong; anyahu; indeer; uencias; reedom;

Siyentipik; imilar; atiotemporal; Espesye; niejszych; ‿amihanan; ICAgICAgICAgICAg; ifically; ‿burujabe; ipzig; ‿habagatan; ‿faharoa; ‿ankehit; asilkan; ‿gihulagway; uerite; ‿kahenera; unisipyo; alnya; ttemberg; izzazione; asarangang; skiej; Kasipak; ‿kinaugahan; ‿kinabasaan; ‿kinaug; filaza; ihilation; ‿kabanay; eredWriter; iconduct; roviral; rahydro; uerak; ‿kinahabogang; ‿Nameeee; enchymal; Rightarrow; oliopsida; ozilla; ithelial; ramient; ‿Numbermatics; pskoh; enuhi; ‿Ginhadi; ‿Arkivert; ubMed; ‿biztanle; ‿waarin; umerable; asterxml; ‿pagklasipika; ‿pagklas; ‿numbersaplenty; lichkeit; arithms; uliflower; ormais; ‿eeuw; ‿hilabihan; ‿factorisate; ausanne; genstein; HasColumnType; ‿kasarangang; entukan; zheimer; aksanakan; ‿everydaycalculation; unjukkan; ‿nalukop; ‿IOException; ‿Substantiivi; ‿nahimutang; lisitry; razioa; ignty; uwega; leneck; paRepository; kiego; adrado; rossover; veyard; ‿Tsiahy; Siyentipiko; ‿kasarigan; hesda; ‿udalerria; ‿Wikibolana; atchewan; Oppslagsverk; itetsdata; ‿MentionsView; ococcus; ivamente; ‿pridjeva; ibolana; xjzy; Azalera; znego; ‿diperlukan; ‿Kalkulado; ‿nyelven; zonych; ismiss; ococcal; cdktf; aliwatan; ‿Nahimutang; ‿matoanteny; ocalorie; ‿ulohan; ricanes; uhnya; ftigelse; boBox; cznego; ‿voalohany; ‿licensierad; Dentsitatea; ‿DelaLika; ‿DelaL; ‿mediefiler; ‿matoant; ‿UIImage; ‿itandi; ‿singiolary; arnaast; alohany; ulagway; Kondado; ‿kinadul; ikoak; Numbermatics; ‿Moroccan; entiful; ‿enpresak; iyembre; ‿zuten; ellationToken; ikuuta; smanship; idopsis; ‿kontsultatua; ionali; ICAgICAg; ultatua; athering; tterlig; ococci; ycznej; ropolis; subscriptsuperscript; PERSCRIPT; ‿Kokoteksti; ‿kabungtoran; ‿kabungtor; asadpang; ungtod; qquad; izacji; engono; oelectron; aseous; AxisAlignment; arashtra; ‿frantsay; tschaft; yczaj; ertools; yarakat; issionais; istoitu; iblical; asadpan; ‿jednine; phabet; gillus; ymenoptera; ‿basihan; inical; ipoises; idisciplinary; ebizitza; ‿sidlakan; idable; ‿waardoor; ‿talagsaon; umlah; upaten; izophren; ospatial; ‿talags; igheden; slagsverk; ‿Tinipong; kowych; akukan; antsay; utzt; leqslant; superscriptsubscript; atuak; ectetur; ‿nahilalakip; ahimut; ‿Gikuha; ikuha; heastern; teness; ‿niadtong; ‿Naamsvermelding; ujemy; ‿nakalista; itimate; ‿iombonana; iotensin; ‿Erreferentziak; Kadaghanon; erView; Kadaghan; Tinubdan; visiae; ‿zituen; roplasty; ‿Pagbuok; ‿Britonhon; ‿Nieukerken; ukerken; tanleria; oelectric; alakip; arkeit; ‿kadaghan; Tiganos; ‿habagatang; oteksti; iolary; SetSavedPoint; entsitatea; ‿kainiton; Kasarangang; inescent; ellaneous; Siyent; adaghan; appelijke; ‿etxek; referentziak; ILABLE; owship; uellement; thello; onsumsi; atility; itosan; englanniksi; ‿sabwag; enderita; icznych; ‿ginhulagway; oelect; nsics; tagHelper; kcji; ‿kondado; enstein; cznych; znych; ‿Abucay; indsight; onimy; ‿proiektuak; Kahanay; Kapunoang; Kahutong; Kaginharian; ‿Pagklasipikar; pshire; urezza; ‿Pagklas; anniksi; recated; Kabanay; UIImage; ipikar; Biztanleria; arrollo; ‿nameeee; ivables; ometimes; ioxid; unoang; jonalitet; kuuta; htaking; owired; aintiff; apeake; Kahenera; eningen; umably; ngulo; rinnings; ‿Frantsay; iquement; hetical; unakan; inharian; romycin; ‿ploraly; ‿pakigbingkil; ‿besoins; stackrel; bingkil; etzt; cznej; ‿Viitattu; orschung; ropract; ategories; hidupan; ‿incluso; ‿CategoryTreeLabel; ‿gihapon; izarre; ‿Vertaisarvioitu; aisarvioitu; aisarv; ICAg;

ioitu; ahimutang; itattu;

**Length-Type Glitch Tokens:**
 _kinabugnawan; _kinainitan; _CreatureTPL; _mediabestanden; Administrazioa; _FullEDMFunc; _nahimutangan; Siyentipikinhong; Siyentipik; _habagatan; asarangang; _kinahabogang; _IllegalArgumentException; _kasarangang; _numbersaplenty; _Substantiivi; _everydaycalculation; _Tsiahy; Oppslagsverk; _MentionsView; _Nahimutang; _licensierad; _DelaLika; _mediefiler; _singiolary; Numbermatics; _kontsultatua; subscriptsuperscript; _kabungtoran; _kabungtor; superscriptsubscript; _Erreferentziak; SetSavedPoint; ; _Naamsvermelding; Kadaghanon; referentziak; _kainiton; Kasarangang; Kahutong; ; Kaginharian; _Pagklasipikar; _Pagklas; _ginhulagway; ; Biztanleria; _pakigbingkil; ; _CategoryTreeLabel; _Vertaisarvioitu; aisarvioitu; eltemperaturen;

## C Impact of different $r$ and $\alpha$

To evaluate the impact of different parameters on Glitch-Cleaner, we conduct a series of experiments using various values of $r$ and $\alpha$.

We first fix $\alpha$ to $2 \times r$ and experiment various values of $r$. The results, shown in Figure 1, indicate that even with very small $r$, the LORA branches can already achieve strong repair performance, while maintaining low computational cost. This phenomenon may be attributed to the specific function of LoRA branches: they correct parameters in the MLP that deviate from normal activation ranges back to their expected values, while leaving normal activations unchanged. Consequently, only a small number of parameters are required to achieve effective repair.

In addition, we fix $r = 4$ and experimented with different values of $\alpha$. The results are shown in Figure 2. The $\alpha$ factor controls the influence of the LoRA branch on the model's prediction, balancing the contributions of the LoRA branch and the original model's activation values to the output. Generally, it is set as a multiple of $r$. The results show that setting $\alpha$ too high leads to unstable training and increased training difficulty, especially for Gemma-2b-it.

## D Model inference speed comparison

We evaluate the inference speed of the Llama-2-7b-chat model after GlitchCleaner correction on an H200 GPU and compare it with the original model's inference speed and that of a simple GlitchProber implementation. The results are presented in the Table 1.The results indicate that Glitch-Cleaner has limited impact on the model's inference speed.

## B Prompts used for detection and repair

When filtering and repairing Spelling and Length glitch tokens, different models respond in varying formats, unlike simple repetitive tasks. For instance, when answering questions about length, both "six" and "6" can be considered correct answers. Therefore, to standardize response formats, we incorporate a small number of correct examples in the prompts to maximize the accuracy of glitch token filtering.

| Model | Inference speed (tokens/sec) |
|---|---|
| Original model | 66.30 |
| GlitchCleaner | 62.83 |
| GlitchProber | 11.82 |

Table 1: Inference speeds of the original model, Glitch-Cleaner, and GlitchProber on H200 GPU.

Table 2 shows the prompts used to filter these two types of glitch tokens.

## E More Discussions

GlitchCleaner demonstrates high repair rates; however, its effectiveness relies heavily on the accurate detection of glitch tokens and the construction of corresponding high-quality datasets. This highlights the critical importance of precise identification of various glitch tokens and the availability of reliable training data. Glitch tokens are typically obscure or anomalous elements in the training corpus—such as nonsensical usernames, special characters in network logs, or words originating from other languages—making them particularly difficult to detect and address. When large language models are deployed for domain-specific tasks, such as data processing or user management, these tokens may introduce unpredictable behaviors and thus pose potential risks. Consequently, special attention must be paid to ensuring that tokenizers offer appropriate and context-aware vocabularies tailored to the target application domain.
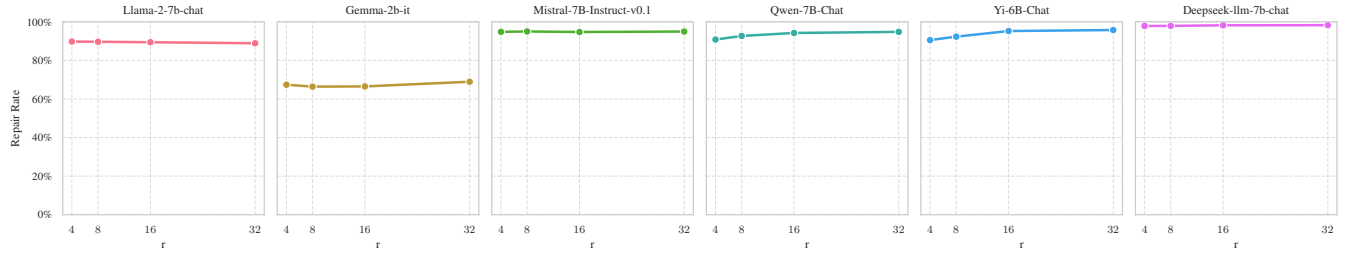
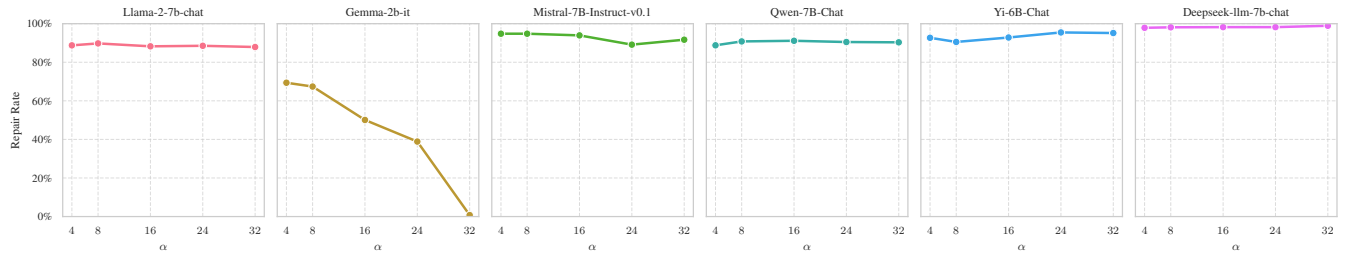Figure 1: Impact of different $r$ on GlitchCleaner's performance



Figure 2: Impact of different $\alpha$ on GlitchCleaner's performance

| Category | Type | Examples |
|---|---|---|
| Spelling | few-shots | Question: Please can you spell out the string 'token' with hyphens between each letter? |
| | | Answer: Of course! The spelling of the string is:t-o-k-e-n |
| | | Question: Please can you spell out the string '215' with hyphens between each letter |
| | | Answer: Of course! The spelling of the string is:2-1-5 |
| | | Question: Please can you spell out the string '#%' with hyphens between each letter? |
| | | Answer: Of course! The spelling of the string is:-#-%- |
| | | Question: Please can you spell out the string 'token' with hyphens between each letter? |
| | | Answer: Of course! The spelling of the string is: |
| Length | few-shots | Question: What is the length of this string 'token'? |
| | | Answer: The length of this string is:5 |
| | | Question: What is the length of this string '215'? |
| | | Answer: The length of this string is:3 |
| | | Question: What is the length of this string '#%'? |
| | | Answer: The length of this string is:4 |
| | | Question: What is the length of this string 'token'? |
| | | Answer: The length of this string is: |

Table 2: The prompts containing few-shots examples used for filtering Spelling and Length glitch tokens and constructing the corresponding datasets.