

Tokens with Meaning: A Hybrid Tokenization Approach for NLP

Anonymous ACL submission

001

Abstract

002 Tokenization plays a pivotal role in natural
003 language processing (NLP), shaping how
004 textual data is segmented, interpreted, and
005 processed by language models. Despite the
006 success of subword-based tokenization tech-
007 niques such as Byte Pair Encoding (BPE)
008 and WordPiece, these methods often fall
009 short in morphologically rich and aggluti-
010 native languages due to their reliance on
011 statistical frequency rather than linguistic
012 structure. This paper introduces a linguis-
013 tically informed hybrid tokenization frame-
014 work that integrates rule-based morpholog-
015 ical analysis with statistical subword seg-
016 mentation to address these limitations. The
017 proposed approach leverages phonological
018 normalization, root-affix dictionaries, and a
019 novel tokenization algorithm that balances
020 morpheme preservation with vocabulary ef-
021 ficiency. The framework also incorporates
022 special tokens for whitespace and ortho-
023 graphic case, including an <uppercase>
024 token to prevent vocabulary inflation from
025 capitalization. Byte Pair Encoding is inte-
026 grated to support out-of-vocabulary cover-
027 age without compromising morphological
028 coherence. Evaluation on the TR-MMLU
029 benchmark—a large-scale, Turkish-specific
030 NLP benchmark—demonstrates that the
031 proposed tokenizer achieves the highest
032 Turkish Token Percentage (90.29%) and
033 Pure Token Percentage (85.8%) among
034 all tested models. Comparative analysis
035 against widely used tokenizers from models
036 such as LLaMA, Gemma, and OpenAI’s
037 GPT reveals that the proposed method
038 yields more linguistically meaningful and
039 semantically coherent tokens. A qualitative
040 case study further illustrates improved mor-
041 pheme segmentation and interpretability
042 in complex Turkish sentences. This work
043 contributes to ongoing efforts to improve
044 tokenizer design through linguistic align-
045 ment, offering a practical and extensible

solution for enhancing both interpretabil-
ity and performance in multilingual NLP
systems.

Keywords: Tokenization, Morphologically
Rich Languages, Morphological Segmenta-
tion, Byte Pair Encoding, Turkish NLP,
Linguistic Integrity, Low-Resource Lan-
guages

1 Introduction

Tokenization is a foundational step in Natural
Language Processing (NLP), directly impacting
vocabulary construction, model efficiency, and
downstream task performance (Liu et al., 2019).
While subword-based methods like Byte Pair
Encoding (BPE) (Sennrich et al., 2016) and
WordPiece (Schuster and Nakajima, 2012) effec-
tively handle out-of-vocabulary (OOV) words
in high-resource languages, they often disregard
the linguistic structure of morphologically rich
languages. In agglutinative languages like Tur-
kish, Finnish, and Hungarian, words are formed
by appending multiple affixes to a root, result-
ing in complex surface forms. Frequency-based
subword models frequently violate morphemic
boundaries in these languages, reducing seman-
tic coherence and interpretability (Toraman
et al., 2023).

Turkish poses specific challenges due to its
agglutinative nature and phonological processes
such as vowel harmony and consonant alterna-
tion. Words are formed by appending multi-
ple affixes to a root, producing an expansive
set of surface forms. For instance, the single
word *Avrupalılaştıramadıklarımızdanmışsınızcasına* ("as if you were one of those whom we
could not make resemble a European") conveys
a meaning that requires an entire sentence in
English. Standard tokenizers often treat these
variants as distinct or fragment them inconsis-
tently, leading to vocabulary redundancy and

046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085

poor alignment with linguistic units (Bayram et al., 2025b). Recent benchmarks, such as TR-MMLU (Bayram et al., 2025a), indicate that "token purity"—the alignment of tokens with morphemes—correlates strongly with model performance. Token purity fundamentally determines the clarity of the statistical patterns LLMs learn. Unlike impure subwords that introduce ambiguity, tokens aligned with complete morphemes provide consistent semantic and syntactic signals, facilitating better generalization across complex word forms (Hofmann et al., 2021). Empirical evidence supports this: morphologically informed models have been shown to outperform standard BPE baselines, achieving superior efficiency even with fewer training iterations (Jabbar, 2024). This principle mirrors object-centric learning in computer vision, where decomposing inputs into meaningful entities—rather than undifferentiated features—enhances recognition.(capsule networks (Sabour et al., 2017) and object-centric architectures like Slot Attention (Locatello et al., 2020)) Consequently, token purity is not merely a linguistic preference but a structural necessity for semantic awareness, motivating its use as a primary evaluation metric in this work.

To address these limitations, we introduce a linguistically informed hybrid tokenization framework. Our approach integrates rule-based morphological segmentation with BPE to ensure both linguistic fidelity and broad coverage. First, phonological processes unify surface variants into shared identifiers. Second, a dedicated <uppercase> token decouples capitalization from lexical identity, preventing vocabulary inflation. Third, explicit formatting tokens are employed to preserve the structural integrity of the input for layout-sensitive tasks. Finally, a hybrid architecture integrates dictionary-based morphological segmentation with BPE, balancing the need for linguistic purity with robust coverage for out-of-vocabulary terms.

2 Related Work

Tokenization significantly impacts model performance, especially in morphologically rich languages (Toraman et al., 2023). While subword methods like BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012) are standard, they often fail to capture the ag-

glutinative structure of languages like Turkish, Finnish, and Hungarian (Baykara and Güngör, 2022).

Early Turkish NLP relied on rule-based morphological analyzers like Zemberek (Akin and Akin, 2007), which offered precise segmentation but lacked the scalability required for modern LLMs. Recent research has sought to bridge this gap. Toraman et al. (2023) showed that morphological tokenization could recover 97% of BERT’s performance with a fraction of the model size. Similarly, Pan et al. (2020) and Huck et al. (2017) demonstrated that morphology-aware segmentation improves Neural Machine Translation (NMT) by reducing data sparsity. Most recently, (Asgari et al., 2025) introduced MorphBPE, a hybrid method that constrains BPE merges to respect morpheme boundaries. Their experiments demonstrated that this linguistic alignment yields tangible computational benefits, resulting in significantly lower training loss and faster convergence rates compared to standard subword models.

Hybrid approaches combining rule-based and statistical methods have shown promise. Kayali and Omurca (2024) used a hybrid tokenizer for Turkish NER and summarization, finding benefits in preserving linguistic structure. Jabbar (2024) introduced MorphPiece for English, achieving superior performance with smaller vocabularies. However, challenges remain in balancing vocabulary size, sequence length, and computational efficiency (Henderson et al., 2022; Kaya and Tantug, 2024). Our work extends these efforts by introducing a fully hybrid pipeline that integrates phonological normalization directly into the tokenization process.

Despite the theoretical appeal of morphological segmentation, some studies argue that strict linguistic boundaries may not always be optimal for neural models. Kudo (2018) introduced "subword regularization," suggesting that exposing models to multiple segmentations of the same word (including non-canonical ones) can improve robustness and generalization. Similarly, recent large-scale multilingual models often favor larger, data-driven vocabularies that maximize compression rates over linguistic purity. This "byte-premium" hypothesis suggests that minimizing the number of tokens per word is the primary driver of cross-linguistic perfor-

188 mance gaps (Martins et al., 2024).
189

190 Furthermore, the rise of massive multilingual
191 LLMs presents a dilemma for language-specific
192 tokenization. While a dedicated Turkish to-
193 kenizer offers superior alignment for Turkish
194 text, it may not be easily integrated into a
195 model trained on 100+ languages without sig-
196 nificantly increasing the combined vocabulary
197 size or complicating the embedding space. Our
198 work acknowledges this tension but argues that
199 for high-stakes or specialized applications in
200 morphologically rich languages, the benefits
201 of semantic coherence and interpretability out-
202 weigh the costs of language-specific engineering.

203 3 Methodology

204 Traditional subword tokenization methods like
205 BPE (Sennrich et al., 2016) and WordPiece
206 (Schuster and Nakajima, 2012) often fail to
207 capture the rich internal structure of agglu-
208 tinative languages. For example, the Turkish
209 word *anlayabildiklerimizden* is composed of mul-
210 tiple morphemes (*anla-yabil-dik-ler-imiz-den*).
211 Standard tokenizers fragment such words ar-
212 bitrarily, obscuring grammatical function. To
213 address this, we propose a hybrid tokenization
214 framework that prioritizes morphological seg-
215 mentation while retaining BPE as a fallback
for robustness.

216 3.1 Dictionary Construction

217 The foundation of our hybrid tokenizer is a cu-
218 rated set of morphological dictionaries derived
219 from open-source linguistic resources, primarily
220 the Zemberek NLP framework (Akin and Akin,
221 2007) and the Turkish Language Association
222 (TDK) data.

223 1. **Root Dictionary** : We extracted ap-
224 proximately 22,000 high-frequency Turk-
225 ish roots (nouns and verbs) from a large-
226 scale corpus of Turkish web text. These
227 roots were filtered to exclude rare or ar-
228 chaic terms that would unnecessarily in-
229 flate the vocabulary. Special control tokens
230 (<uppercase>, <unknown>, <pad>, <eos>)
231 were added to support model training.

232 2. **Suffix Dictionary** : This dictionary
233 contains 230 distinct derivational and in-
234 flectional suffixes. Crucially, we applied
235 phonological abstraction: suffixes that are

236 phonetically distinct but functionally iden-
237 tical (allomorphs) are mapped to a single
238 canonical ID. For example, the plural suffix
239 forms *-lar* and *-ler* share the same token
240 ID, as do the four variants of the accusative
241 case (*-i*, *-i*, *-u*, *-ü*). Similarly, the future
242 tense suffixes *-acak* and *-ecek* are unified
243 under a single canonical ID, reflecting their
244 shared grammatical function despite the
245 consonant change.

246 3. **BPE Vocabulary** : To ensure full cover-
247 age for foreign words, proper names, and
248 rare scientific terms, we trained a Byte
249 Pair Encoding (BPE) model on the same
250 corpus with a vocabulary limit of 10,000
251 subwords. This serves as a fallback mech-
252 anism for any segment not found in the
253 root or suffix dictionaries.

254 3.2 Encoding Process

255 Encoding process consists of three main stages
256 when converting text into sequences of IDs.
257 First, the input text is separated into base
258 words by splitting it based on space char-
259 acters. In this stage, a space prefix is added to
260 the beginning of each word (Example: "elma"
261 → " elma"), due to the fact that all tokens
262 in the tokenizer's roots dictionary include a
263 leading space. Second, these separated words
264 are processed according to their capitalization
265 status. If a word starts with a capital letter,
266 it is represented by a special <uppercase>
267 token followed by the lowercase version of the
268 word (Example: "Ankara" → <uppercase> +
269 " ankara"). If capitalization is not at the begin-
270 ning (like "iPhone"), the word is fragmented
271 and the capitalized parts are parsed with spe-
272 cial tokens. The final stage is converting the
273 pieces into numerical IDs. All obtained pieces
274 (roots, suffixes, and BPE tokens) are searched
275 within the root, suffix, and BPE dictionaries.
276 The search algorithm proceeds by iteratively
277 removing letters backward from the end of the
278 piece (using the longest prefix match logic). For
279 instance, in the word "kitaplar," attempts are
280 made sequentially with "kitaplar," "kitapla,"
281 etc.; once the root "kitap" is found in the root
282 list, its ID is added, and then the remaining
283 part, "lar," is searched in the suffix list and its
284 token ID is added. Upon completion of these
285 operations, the final numerical ID list of the

286 text is generated.

318

```
Algorithm: Hybrid Encoding
Input: text string T
Output: list of token IDs
```

319

1. Split T into words by whitespace.
2. For each word W:
 - a. Identify uppercase positions.
 - b. Split W into segments (handle camelCase).
 - c. For each segment S:
 - i. Check ROOTS for longest prefix match.
 If match: add ID, continue.
 - ii. Check SUFFIXES for longest prefix match.
 If match: add ID, continue.
 - iii. Check BPE for longest prefix match.
 If match: add ID, continue.
 - iv. Else: add <unknown> ID.
 - d. Insert <uppercase> tokens based on (a).
3. Return list of IDs.

287 Figure 1: Pseudocode for the hybrid encoding process.
288289

3.3 Decoding and Phonological 290 Phonological Resolution

291 The **Decoding** phase reconstructs text from
292 token IDs using a "Single ID, Multiple Views"
293 principle. Since multiple surface forms (allo-
294 morphs) map to a single ID, the decoder must
295 dynamically resolve the correct form based on
296 context. This process is critical for generating
297 natural-sounding Turkish text.298

3.3.1 Root Resolution (Lookahead)

299 In this system, a special root selection mech-
300 anism is employed in the Decoder stage, fol-
301 lowing tokenization, to accurately model mor-
302 phological changes in natural language. Due to
303 the morphological structure of Turkish, surface
304 forms of the same root that have undergone
305 sound events (phonological changes) are rep-
306 resented by the same ID in the system, even
307 though they look different. For instance, kitap
308 (book) and kitab- (the form resulting from
309 consonant softening), ağız (mouth) and ağız- (the
310 form resulting from vowel deletion), and küçük
311 (small) and küçüğ- (the form resulting from con-
312 sonant deletion/softening) all correspond to the
313 same root ID . As seen in these examples, even
314 though the forms affected by sound events are
315 different on the surface, the model unifies these
316 variants under a single root identity. While the
317 Tokenizer stores these potential variations as a
list associated with the ID, the Decoder's main
task is to select the correct surface form from318 this list that is most appropriate for the con-
319 text. To make this selection, the Decoder calls
320 the specially defined function for every root ID.
321 This function determines the most correct mor-
322 phological form by looking at the information
323 of the suffix that will follow the root. The main
324 factors the function considers are: whether the
325 next token is a suffix, if it is a suffix, whether
326 it starts with a vowel (like "yor," "acak," or
327 "1"), whether it is one of certain specific special
328 suffixes, and finally, whether the root's ID falls
329 within a range subject to special rules.330 A concrete example of this root selection
331 mechanism is as follows: Let's assume the
332 tokens corresponding to Root ID 100 are →
333 ["sicak", "sicağ", "sica"]. If the ID sequence
334 arrives as [100, 2034] (where ID 2034 repre-
335 sents the suffix "1"), the Decoder immediately
336 activates the `_select_correct_root` function.
337 This function detects that ID 2034 is a vowel-
338 initial suffix. Based on this contextual infor-
339 mation, the decision is made that consonant
340 softening must occur, and the softened form of
341 the root, "sicağ," is selected. When the suffix is
342 subsequently added ("sicağ" + "1"), the correct
343 morphological form, "sıcığı" (the accusative
344 case of hot/warm), is obtained. This ensures
345 the root is selected and prepared in the correct
346 form based on the type of suffix that will follow
347 it. Additionally, the Decoder manages the Up-
348 percase Marker, which is a special grammatical
349 token. When the token ID is 0, the Decoder
350 understands that this marker must convert the
351 first letter of the immediately following root
352 into a capital letter.353

3.3.2 Suffix Resolution (Lookbehind)

354 In this system, if a token ID is 20,000 or greater,
355 it represents a suffix, and just like roots, any
356 given suffix ID may correspond to multiple sur-
357 face forms. This is due to the fact that Tur-
358 kish suffixes follow phonological harmony rules
359 such as major/minor vowel harmony and con-
360 sonant alternation (e.g., the plural suffix 20000 →
361 ["lar", "ler"], or the future tense suffix 20030 →
362 ["acak", "ecek", "acağ", "eceğ", "yacak", "yecek",
363 "yacağ", "yeceğ"]). The primary responsibility
364 of the Decoder is to call a specialized function
365 that selects the correct surface form from these
366 alternative suffix lists based on context. This
367 function evaluates several factors when choos-
368 ing the appropriate suffix: the last vowel of

the preceding word to ensure vowel harmony, whether the final letter of the preceding word is a voiceless consonant for consonant alternation, whether the next token begins with a vowel to determine if a buffer consonant (such as *y*) is required, the position of the suffix within the word, and whether the suffix belongs to a special category (e.g., “la/le,” “da/de-ta/te,” “cik/cik,” “mak/mek,” “acak/ecek,” etc.). After retrieving the list of surface variants associated with the suffix ID, the function routes the selection process to the appropriate specialized sub-function depending on the suffix category, ensuring that the correct surface form is chosen.

- **Example 1 — Vowel Harmony:** If the preceding word is “ev” (its last vowel “e” is a front vowel), the plural suffix (20000) selects the form “ler” to satisfy front-vowel harmony. Result: [“ev”, “20000”] → “evler”.
- **Example 2 — Consonant Assimilation (Fortition):** If the preceding word is “kitap” (ending with the voiceless consonant “p”) and the suffix is -da/-de, the Decoder applies consonant assimilation and hardens the suffix to “ta/te”. Since the last vowel of the word is “a” (a back vowel), the suffix is further matched to the back-vowel form. Result: “kitapta”.
- **Example 3 — Future Tense (Vowel Harmony):** If the preceding word is “bak” (its last vowel “a” is a back vowel) and the suffix is -acak/-ecek, the Decoder selects the form “acak” to satisfy back-vowel harmony. Result: “bakacak”.

4 Results and Analysis

We evaluated the proposed tokenizer on the TR-MMLU benchmark (Bayram et al., 2025a), comparing it against five state-of-the-art tokenizers: `google/gemma-2-9b`, `meta-llama/Llama-3.2-3B`, `Qwen/Qwen2.5-7B-Instruct`, `CohereForAI/aya-expanse-8b`, and `microsoft/phi-4`. Metrics include Vocabulary Size, Total Tokens, Unique Tokens, Turkish Token Percentage (TR %), and Pure Token Percentage (Pure %).

As shown in Table 1, our tokenizer achieves the highest linguistic alignment (90.29% TR %, 85.80% Pure %) while being 8× smaller than

Table 365: Performance comparison on the TR-MMLU dataset.
370

Tokenizer	Vocab	Tokens	Unique	TR %	Pure %
<code>turkish_tokenizer</code> (Ours)	32k	707k	11.1k	90.29	85.80
<code>google/gemma-2-9b</code>	256k	497k	6.3k	40.96	28.49
<code>meta-llama/Llama-3.1</code>	128k	488k	6.8k	45.77	31.45
<code>Qwen/Qwen2.5-7B-Instruct</code>	152k	561k	5.7k	40.39	30.15
<code>CohereForAI/aya-expanse-8b</code>	256k	434k	8.5k	53.48	32.96

Gemma’s 256k tokenizer. This size efficiency is crucial, as it significantly reduces the embedding layer parameters, thereby minimizing the model’s memory footprint. This linguistic precision leads to an anticipated trade-off: an increased sequence length, with our tokenizer generating 63% more tokens `aya-expanse` (707k vs. 434k). While a higher token count means longer sequences and a higher computational cost for attention mechanisms, it provides a crucial benefit for agglutinative languages. It furnishes the model with a more explicit and regular representation of the language, allowing it to learn the compositional rules of the grammar rather than being forced to memorize millions of distinct surface forms (e.g., *geldim*, *geldin*, *geldi...*) as individual vocabulary items.

These results confirm that frequency-based tokenizers, even with large vocabularies, struggle with the agglutinative structure of Turkish. Our hybrid approach, by contrast, yields tokens that are linguistically meaningful and consistent, a success that underscores its value for managing this characteristic complexity. By ensuring our tokens are linguistically meaningful and respect morpheme boundaries, we not only achieve higher alignment metrics but also provide the model with a superior, more regular representation of the language. This allows the model to efficiently learn the compositional nature of Turkish grammar, demonstrating that prioritizing morphological integrity can be more effective for complex linguistic structures than simply relying on large, frequency-driven vocabularies.

5 Discussion

The current practice in multilingual models (e.g., 4LLaMA-3, Gemma-2) is to use large, shared vocabularies (128k–256k), which, while efficient for high-resource languages like English, often treats low-resource, agglutinative languages as “second-class citizens” by forcing them into fragmented subword sequences.

Our findings propose a viable alternative: a modular tokenization approach or a "mixture-of-tokenizers" architecture, where language-specific morphological adapters are swapped in during training . This method allows models to process each language in its most natural structural form. Crucially, the principles of our hybrid framework are highly portable; they are directly applicable to other agglutinative languages such as Finnish, Hungarian, and Estonian, which share features like rich suffixation and compounding. The underlying "Longest Prefix Match" algorithm with phonological abstraction is language-agnostic, requiring only the replacement of language-specific dictionaries and the definition of phonological rules (like vowel harmony groups). This inherent portability is a significant advantage over purely statistical methods, which demand massive corpus retraining to implicitly discover these linguistic rules.

6 Conclusion

In this study, we introduced a linguistically informed hybrid tokenization framework designed to address the challenges of morphologically rich languages. By integrating rule-based morphological analysis with BPE, our approach preserves morpheme boundaries and minimizes vocabulary redundancy. Empirical evaluations on the TR-MMLU dataset demonstrate that our tokenizer achieves significantly higher linguistic alignment (90.29% TR %, 85.80% Pure %) compared to state-of-the-art multilingual models like LLaMA, Gemma, and Qwen. These results validate that incorporating linguistic structure into tokenization yields more semantically coherent representations.

6.1 Limitations

While our results are promising, this study has several limitations. First, our evaluation relies primarily on intrinsic metrics (TR % and Pure %). Due to computational constraints, we did not pretrain a language model from scratch to empirically verify the impact of our tokenizer on downstream task performance (e.g., perplexity, classification accuracy). Establishing a causal link between token purity and model performance remains a critical next step. Second, our current implementation is in Python, which may not match the inference speed of

highly optimized Rust-based tokenizers used in production LLMs. We have not yet conducted rigorous benchmarking of processing time or memory usage. Third, the approach relies on manually curated dictionaries, which may require maintenance to cover evolving language use and neologisms.

6.2 Future Work

Future research will focus on three key areas: (1) **Downstream Evaluation:** Training a small-scale language model (e.g., 100M parameters) using our tokenizer to measure improvements in perplexity and task-specific accuracy compared to standard BPE; (2) **Optimization:** Re-implementing the tokenizer in Rust to ensure it meets the latency requirements of real-time applications; and (3) **Generalization:** Extending the framework to other agglutinative languages such as Finnish and Hungarian to test the cross-linguistic validity of our hybrid approach.

References

- Mehmet Dündar Akin and Ahmet Afşin Akin. 2007. Türk Dilleri İçin Açık Kaynaklı Doğal Dil İşleme Kütüphanesi: ZEMBEREK. *elektrik mühendisliği*, 431:38–44.
- Ehsaneddin Asgari, Yassine El Kheir, and Mohammad Ali Sadraei Javaheri. 2025. Morphbpe: A morpho-aware tokenizer bridging linguistic complexity for efficient llm training across morphologies. *arXiv preprint*. ArXiv:2502.00894 [cs.CL].
- Batuhan Baykara and Tunga Güngör. 2022. Abstractive text summarization and new large-scale datasets for agglutinative languages turkish and hungarian. *Language Resources and Evaluation*, 56(3):973–1007.
- M. Ali Bayram, Ali Arda Fincan, Ahmet Semih Gümüş, Banu Diri, Savaş Yıldırım, and Öner Aytaş. 2025a. Setting standards in turkish nlp: Tr-mmlu for large language model evaluation. *arXiv preprint*. ArXiv:2501.00593 [cs].
- M. Ali Bayram, Ali Arda Fincan, Ahmet Semih Gümüş, Sercan Karakas, Banu Diri, and Savaş Yıldırım. 2025b. Tokenization standards for linguistic integrity: Turkish as a benchmark. *arXiv preprint*. ArXiv:2502.07057 [cs].
- Peter Henderson, Jieru Hu, Joshua Romoff, Emma Brunskill, Dan Jurafsky, and Joelle Pineau. 2022. Towards the systematic reporting of the energy and carbon footprints of machine learning. *arXiv preprint*. ArXiv:2002.05651 [cs].

- 562 Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. 2021. Superbizarre is not superb:
563 Derivational morphology improves bert’s interpretation of complex words. In *Proceedings of*
564 *the 59th Annual Meeting of the Association for*
565 *Computational Linguistics*, pages 3594–3608.
566
567
- 568 Matthias Huck, Simon Riess, and Alexander Fraser.
569 2017. Target-side word segmentation strategies
570 for neural machine translation. In *Proceedings of*
571 *the Second Conference on Machine Translation*,
572 pages 56–67.
- 573 Haris Jabbar. 2024. Morphpiece: A linguistic tok-
574 enizer for large language models. *arXiv preprint*.
575 ArXiv:2307.07262 [cs].
- 576 Yigit Bekir Kaya and A. Cüneyd Tantug. 2024. Ef-
577 fect of tokenization granularity for turkish large
578 language models. *Intelligent Systems with Appli-*
579 *cations*, 21:200335.
- 580 Nihal Zuhal Kayali and Sevinç İlhan Omurca. 2024.
581 Hybrid tokenization strategy for turkish abstrac-
582 tive text summarization. In *2024 8th Interna-*
583 *tional Artificial Intelligence and Data Processing*
584 *Symposium (IDAP)*, pages 1–6.
- 585 Taku Kudo. 2018. Subword regularization: Im-
586 proving neural network translation models with
587 multiple subword candidates. *arXiv preprint*.
588 ArXiv:1804.10959 [cs].
- 589 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei
590 Du, Mandar Joshi, Danqi Chen, Omer Levy,
591 Mike Lewis, Luke Zettlemoyer, and Veselin
592 Stoyanov. 2019. Roberta: A robustly opti-
593 mized bert pretraining approach. *arXiv preprint*.
594 ArXiv:1907.11692 [cs].
- 595 Francesco Locatello, Dirk Weissenborn, Thomas Un-
596 terthiner, Aravindh Mahendran, Georg Heigold,
597 Jakob Uszkoreit, Alexey Dosovitskiy, and
598 Thomas Kipf. 2020. Object-centric learning with
599 slot attention. In *Advances in Neural Infor-*
600 *mation Processing Systems*, volume 33, pages 11525–
601 11538.
- 602 Pedro Henrique Martins, Patrick Fernandes,
603 João Alves, Nuno M. Guerreiro, Ricardo Rei,
604 Duarte M. Alves, José Pombal, Amin Farajian,
605 Manuel Faysse, Mateusz Klimaszewski, Pierre
606 Colombo, Barry Haddow, José G. C. de Souza,
607 Alexandra Birch, and André F. T. Martins. 2024.
608 Eurollm: Multilingual language models for eu-
609 rope. *arXiv preprint*. ArXiv:2409.16235 [cs].
- 610 Yirong Pan, Xiao Li, Yating Yang, and Rui Dong.
611 2020. Morphological word segmentation on agglu-
612 tutive languages for neural machine translation.
613 *arXiv preprint*. ArXiv:2001.01589 [cs].
- 614 Sara Sabour, Nicholas Frosst, and Geoffrey E. Hin-
615 ton. 2017. Dynamic routing between capsules.
616 *arXiv preprint*. ArXiv:1710.09829 [cs].
- 617 Mike Schuster and Kaisuke Nakajima. 2012.
618 Japanese and korean voice search. In *2012 IEEE*
619 *International Conference on Acoustics, Speech*
620 *and Signal Processing (ICASSP)*, pages 5149–
621 5152.
- 622 Rico Sennrich, Barry Haddow, and Alexandra
623 Birch. 2016. Neural machine translation of
624 rare words with subword units. *arXiv preprint*.
625 ArXiv:1508.07909 [cs].
- 626 Cagri Toraman, Eyup Halit Yilmaz, Furkan
627 Sahinuc, and Oguzhan Ozcelik. 2023. Impact
628 of tokenization on language models: An analysis
629 for turkish. *ACM Transactions on Asian and*
630 *Low-Resource Language Information Processing*,
631 22(4):1–21. ArXiv:2204.08832 [cs].