SindBERT, the Sailor: Charting the Seas of Turkish NLP

Raphael Scheible-Schmitt^{1,2} and Stefan Schweter³

¹School of Computation, Information and Technology, Technical University of Munich, Germany, ²Institute of General Practice, Faculty of Medicine and Medical Center, University of Freiburg, Germany, ³Independent Researcher, Holzkirchen, Germany

Correspondence: raphael.schmitt@uniklinik-freiburg.de

Abstract

Transformer models have revolutionized NLP, yet many morphologically rich languages remain underrepresented in large-scale pretraining efforts. With SindBERT, we set out to chart the seas of Turkish NLP, providing the first large-scale RoBERTa-based encoder for Turkish. Trained from scratch on 312 GB of Turkish text (mC4, OSCAR23, Wikipedia), SindBERT is released in both base and large configurations, representing the first large-scale encoder-only language model available for Turkish. We evaluate SindBERT on partof-speech tagging, named entity recognition, offensive language detection, and the TUR-BLIMP linguistic acceptability benchmark. Our results show that SindBERT performs competitively with existing Turkish and multilingual models, with the large variant achieving the best scores in two of four tasks but showing no consistent scaling advantage overall. This flat scaling trend, also observed for XLM-R and EuroBERT, suggests that current Turkish benchmarks may already be saturated. At the same time, comparisons with smaller but more curated models such as BERTurk highlight that corpus quality and diversity can outweigh sheer data volume. Taken together, SindBERT contributes both as an openly released resource for Turkish NLP and as an empirical case study on the limits of scaling and the central role of corpus composition in morphologically rich languages. The SindBERT models are released under the MIT license and made available in both fairseq and Huggingface formats.

1 Introduction

The advent of transformer-based models such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) has reshaped natural language processing (NLP), providing contextualized word representations that generalize across a wide range of tasks. While early efforts focused on English and multilingual approaches, research has consistently

shown that monolingual pre-training on large, high-quality corpora yields superior results for the target language (Delobelle et al., 2020a; Scheible et al., 2024; Scheible-Schmitt and Frei, 2025).

For Turkish NLP, several transformer-based encoders have been introduced in recent years. Notable examples include BERTurk (Schweter, 2025), trained on a 35 GB corpus of Turkish OSCAR, Wikipedia, and OPUS data; ELECTRA (Clark et al., 2020) and ConvBERT (Jiang et al., 2021) models trained on both OSCAR and mC4 (35–242 GB)(Jiao et al., 2020). While these models provide important milestones, most are relatively small encoder models trained with earlier-generation methods or focus on architectures other than RoBERTa. The only RoBERTa models out there were not computed in its fullest extend, but rather with small batch size for relatively small period (Toraman et al., 2023; Tas, 2024). Futher, Turkish still lacks a large-scale, high-quality encoder-only model.

To address this gap, we introduce SindBERT, a RoBERTa-based encoder model pre-trained specifically for Turkish. SindBERT builds on the design principles of the German model Gott-BERT (Scheible et al., 2024) and adapts them to the morphological richness and agglutinative structure of Turkish. We construct a byte-level BPE vocabulary optimized for Turkish, train both base and large variants with fairseq (Ott et al., 2019), and leverage TPUv4 hardware (Jouppi et al., 2023) for efficient large-scale pre-training. SindBERT is designed to combine scalability and reproducibility while directly targeting Turkish, resulting in the first large-scale RoBERTa-style encoder model for Turkish. Our contributions are as follows:

- We release SindBERT_{base} and SindBERT_{large}, trained from scratch on Turkish web-text.
- We benchmark SindBERT against existing Turkish and multilingual models.

2 Related Work

The introduction of transformer-based language models such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) marked a paradigm shift in NLP, enabling significant improvements across a wide range of tasks. Building on these foundations, multilingual extensions such as mBERT and in particular XLM-RoBERTa (Chan, 2020) became widely used as strong generalpurpose baselines across more than 100 languages. At the same time, a wave of monolingual adaptations demonstrated that language-specific pretraining often outperforms multilingual alternatives when sufficient high-quality data is available (Delobelle et al., 2020b; Martin et al., 2020; Chan et al., 2020; Scheible et al., 2024; Scheible-Schmitt and Frei, 2025).

Recently, multilingual encoder-only models have seen a revival. EuroBERT (Boizard et al., 2025) revisits the encoder paradigm with innovations from decoder-only models, introducing a family of multilingual encoders for European and global languages with native support for sequences up to 8,192 tokens. Similarly, mmBERT (Marone et al., 2025) scales encoder pretraining to 3T tokens across 1,800+ languages, introducing novel sampling schedules and showing strong performance on both high- and low-resource languages. These developments highlight that encoder-based architectures remain competitive even in an era dominated by large decoder models.

For Turkish, the first widely adopted transformer encoder was BERTurk (Schweter, 2020), trained on a 35 GB mixture of OSCAR, Wikipedia, OPUS, and additional resources. Variants included cased/uncased models and vocabularies of 32k or 128k tokens. Distilled versions (DistilBERTurk) (Jiao et al., 2020) and subsequent models such as ELECTRA (Clark et al., 2020) and ConvBERTurk expanded the model zoo, with some trained on the Turkish portion of mC4 (up to 242 GB) (Schweter, 2025). These provided important baselines but generally followed smaller encoder configurations or explored alternative pretraining architectures rather than scaling RoBERTa.

Building on this line of work, RoBERTurk (Tas, 2024) introduced a RoBERTa-style encoder specifically adapted for Turkish, showing that refined pre-training objectives and tokenizer design can yield competitive results. In parallel, research has underscored the critical role of tokenization

in morphologically rich languages. Toraman et al. (2023) systematically analyzed the impact of vocabulary size and segmentation strategy, showing that larger vocabularies can notably improve performance in morphosyntactic evaluations. However, all these RoBERTa-based models were not extensively trained, typically using moderate batch sizes and relatively few update steps, resulting in comparatively shallow pretraining regimes.

Taken together, these contributions highlight steady progress in Turkish NLP. However, despite the availability of increasingly large corpora and modern training infrastructure, Turkish has lacked a RoBERTa-based encoder model trained from scratch at scale. SindBERT addresses this gap by providing the first large-scale RoBERTa encoder dedicated to Turkish, trained on modern corpora and released openly to the community.

An overview of existing Turkish transformerbased language models is provided in Table 1.

3 Methods

3.1 Training Data

SindBERT was trained on three Turkish corpora: Wikipedia, OSCAR23 (Jansen et al., 2022), and mC4. The corpus was shuffled and lightly filtered, restricted to the removal of documents containing invalid character encodings. The extracted sizes are approximately 242 GB for mC4, 69 GB for OSCAR, and 0.6 GB for Wikipedia, resulting in a combined pre-training corpus of about 312 GB of Turkish text.

3.2 Pre-processing

Similar to RoBERTa, SindBERT relies on byte pair encoding (BPE) (Radford et al., 2019) for subword segmentation, which directly operates on raw text without the need for pre-tokenization or auxiliary tools such as Moses (Koehn et al., 2007). Since the original GPT-2 tokenizer was designed for English, we instead constructed a tokenizer tailored for Turkish. Following the strategy applied in Gott-BERT (Scheible et al., 2024), we trained a dedicated vocabulary using 40 GB of randomly sampled Turkish text, resulting in a 52k subword inventory optimized for the language. In our experience, sampling around 40 GB of text is already enough for the subword statistics to stabilize, while scaling vocabulary training to the entire corpus would primarily increase computational cost without offering substantial gains. While we did not separately

Model	Architecture	Pre-training Data	Corpus Size
BERTurk _{32k,128k}	BERT base	OSCAR, Wikipedia, OPUS, non-public	35 GB
DistilBERTurk	DistilBERT	Distilled from BERTurk (subset)	7 GB
ELECTRA _{small}	ELECTRA small	OSCAR, Wikipedia, OPUS, non-public	35 GB
ELECTRA _{base}	ELECTRA base	OSCAR, Wikipedia, OPUS, non-public	35 GB
ELECTRA _{mC4}	ELECTRA base	mC4	242 GB
ConvBERTurk	ConvBERT base	OSCAR, Wikipedia, OPUS, non-public	35 GB
ConvBERTurk _{mC4}	ConvBERT base	mC4	242 GB
RoBERTurk	RoBERTa-mid	OSCAR, Turkish C4 subset (1 GB)	28 GB
	(12L, 1024H)		
SindBERT _{base}	RoBERTa base	mC4, OSCAR23, Wikipedia	312 GB
SindBERT _{large}	RoBERTa large	mC4, OSCAR23, Wikipedia	312 GB

Table 1: Overview of models evaluated in this work. We only consider **cased** variants even if uncased versions exist.

evaluate the effect of this adaptation on storage size or downstream accuracy, previous work in Dutch (Delobelle et al., 2020a) and German (Scheible et al., 2024) indicates that language-specific tokenizers can yield improvements in both efficiency and performance.

3.3 Pre-training

Following the setup of GottBERT, we pre-trained both SindBERT_{base} and SindBERT_{large} using the fairseq framework on a 128-core TPUv4 pod (Jouppi et al., 2023). Mixed-precision training (fp16/bfloat16) was not employed, so both models were trained entirely in full precision (fp32). This ensures that training dynamics can be attributed directly to model size, without numerical precision optimizations acting as additional factors.

SindBERT_{base} completed training in approximately 29.2 hours, while SindBERT_{large} required around 6.0 days. We followed the standard RoBERTa pretraining schedule with 100k update steps, a global batch size of 8k, a 10k-step warmup, and polynomial learning rate decay. The base model used a peak learning rate of 0.0004, and the large model 0.00015. Similar to Gott-BERT (Scheible et al., 2024), we evaluated after each epoch and stored checkpoints throughout training. Since the dataset size only permitted roughly four epochs, the final checkpoint coincided with the best-performing one.

3.4 Downstream Tasks

To assess the capabilities of SindBERT, we finetuned the model on a diverse suite of Turkish downstream benchmarks covering sequence labeling, text classification, and linguistic acceptability. Training was performed with the Flair framework (Akbik et al., 2019) v0.15.1, using standardized experiment configurations provided in the repository. Hyperparameter optimization was carried out over batch size and learning rate (Table 2), with training capped at a maximum of 30 epochs and early stopping applied (patience = 3). All models employed a linear learning rate schedule with a 10% warmup phase. We evaluated SindBERT on the following tasks:

Part-of-Speech Tagging We used the concatenation of five Turkish Universal Dependencies (UD) datasets: Atis, BOUN, FrameNet, IMST, and Tourism. This diverse set reflects different domains such as spoken language, newswire, and tourism. Providing a measure of syntactic and morphological coverage, we report model's performance using micro F1.

Named Entity Recognition For NER, we finetuned on the Turkish NER dataset introduced in the WikiANN corpus (Pan et al., 2017) and widely used for multilingual evaluation. We used the splits from Rahimi et al. (2019) and report micro F1 across all entity types.

Offensive Language Detection To evaluate robustness on user-generated content, we employed the OffensEval-TR 2020 dataset (Çöltekin, 2020), a corpus of Turkish tweets annotated for the presence of offensive language. The dataset contains over 31k training and 3.5k test instances, labeled in a binary fashion as either *NOT* (not offensive) or *OFF* (offensive). Mentions and URLs were anonymized during preprocessing (e.g., replaced by @USER or URL), while the tweets otherwise preserve the linguistic and pragmatic properties of social media text. We report performance using macro F1.

Linguistic Acceptability To assess fine-grained grammatical knowledge, we include evaluation on TURBLIMP (Başar et al., 2025), a benchmark of 16 core linguistic phenomena ranging from anaphor agreement and argument structure to scrambling and suspended affixation. Each phenomenon is represented by 1,000 minimal pairs, and models are scored following the BLiMP protocol (Warstadt et al., 2020), i.e., assigning higher probability to the grammatical sentence of each pair. For each model we compute the accuracy within every phenomenon and report the average across all 16 categories as the overall TURBLIMP score. This measure complements PoS tagging, NER, and sentiment classification by probing deeper syntactic and morphosyntactic competence.

3.5 Hyperparameters

We focused our grid search on batch sizes and learning rates, selected based on the most frequent bestperforming values in prior experiments (GottBERT, GeistBERT (Scheible-Schmitt and Frei, 2025); see Table 2). Training was applied to PoS, NER and classification and capped at a maximum of 30 epochs, with early stopping applied using a patience of three epochs. All models employed a linear learning rate schedule with a warmup phase of 10% of the total training steps. All downstream finetuning experiments were conducted with a fixed random seed of 1 for the base models and 42 for the large models. This setup ensures reproducibility and consistency within each scale while maintaining overall comparability across model groups; nonetheless, minor deviations may still arise from seed-related variance (Dodge et al., 2020).

Parameter	Values				
Batch Size	16, 32				
Learning Rate	5e-6, 7e-6, 1e-5, 2e-5, 5e-5				
Epochs	up to 30				
	(Early stopping, patience = 3)				

Table 2: Hyperparameter configurations for downstream fine-tuning. Each model—task combination was trained with all permutations, yielding 10 runs per model and task. Reported scores are averaged across seeds for the best configuration.

3.6 Model Properties

Table 3 summarizes the vocabulary sizes and parameter counts of the Turkish and multilingual models included in our evaluation. The smallest

encoder is ELECTRA_{small} (13.7M parameters), followed by DistilBERTurk (67M). Base-scale Turkish encoders, such as ConvBERTurk (cased and mC4 variants), ELECTRA_{base} (cased and mC4), and BERTurk (cased/uncased), cluster between 106M and 111M parameters with 32k vocabularies. RoBERTurk, another RoBERTa-style encoder with a 50k vocabulary, is slightly larger at 125M parameters. SindBERT_{base} grows further to 126M owing to its 52k vocabulary and extended RoBERTa design.

At the mid-scale, mBERT has 178M parameters with a WordPiece vocabulary of nearly 120k tokens, while the 128k-token BERTurk variants reach 184M. Among larger models, XLM-R_{base} contains 278M parameters, while SindBERT_{large} grows to 357M. The largest encoder considered is XLM-R_{large}, with 560M parameters and a 250k-token vocabulary. All values were extracted using Hugging Face's transformers library.

Table 3: Vocabulary size and total parameter count for Turkish transformer-based models. Values were extracted using Hugging Face's transformers library.

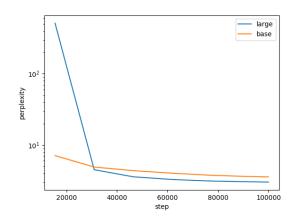
Model	Vocab Size	#Params		
ELECTRA _{small}	32000	13,672,192		
DistilBERTurk	32,000	67,497,984		
ConvBERTurk	32,000	106,815,624		
ConvBERTurk _{mC4}	32,000	106,815,624		
ELECTRA _{base, mC4}	32,000	110,026,752		
BERTurk _{32k}	32,000	110,617,344		
RoBERTurk	50,265	124,644,864		
SindBERT _{base}	52,009	125,985,024		
$mmBERT_{small}$	256,000	140,493,696		
BERTurk _{128k}	128,000	184,345,344		
EuroBERT _{210M}	128,256	211,766,016		
XLM-R _{base}	250,002	278,043,648		
mmBERT _{base}	256,000	306,939,648		
SindBERT _{large}	52,009	357,145,600		
XLM-R _{large}	250,002	559,890,432		
EuroBERT _{610M}	128,256	607,874,688		

4 Results

4.1 Pre-training

During pre-training, we monitored perplexity both on the training set (at each optimization step) and on the validation set (after each epoch; see Figure 1). Across all configurations, the curves follow a consistent convergence pattern. An initial plateau phase can be observed, which is relatively brief for the base models but more pronounced for the large ones. Occasional short upward spikes appear in the training curves; if taken in isolation, these might be misread as divergence, yet they quickly subside as training progresses.

The base models typically stabilize after 20k–30k steps, while the large models require slightly longer but consistently converge by around 40k steps. By the end of training, both configurations achieve comparably low perplexity, underscoring the efficiency of the pre-training setup. This trend is mirrored in the validation perplexity, which shows steady improvements after each epoch. Overall, training perplexity decreased from about 54.5k to 3.93 for the base models and from about 52.2k to 3.24 for the large models, reflecting robust and reliable convergence.



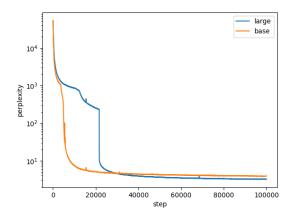


Figure 1: Perplexity of the SindBERT models. Top: validation perplexity measured at checkpoints. Bottom: training perplexity measured at each optimization step.

4.2 Downstream Tasks

Part-of-Speech Tagging Across base-scale models, performance on the Turkish Universal Dependencies treebank is consistently high, with micro-F1 values exceeding 93% for nearly all encoders. The strongest overall results are achieved by ConvBERTurk_{mC4} (94.57), closely followed by SindBERT_{base} (94.47) and BERTurk_{128k} (94.44). Interestingly, both ConvBERTurk variants, trained

with different corpora, maintain a narrow margin over ELECTRA-based and RoBERTa-style encoders, suggesting that architectural innovations like dynamic convolution offer slight but consistent gains in token-level syntactic tagging. The relatively low score of RoBERTurk (87.99) indicates the limitations of early RoBERTa replications for Turkish, likely due to smaller corpora and shorter training schedules. SindBERT_{base} performs competitively within this saturated range, demonstrating strong generalization across tasks despite a larger 52k BPE vocabulary.

Among large-scale encoders, SindBERT_{large} attains the highest F1 (94.63), marginally outperforming XLM-R_{large} (94.39). This indicates that Sind-BERT's pre-training on modern Turkish data contributes positively to syntactic coverage, even when compared to substantially larger multilingual models. The weaker performance of EuroBERT_{610M} (93.33) may reflect its more domain-diverse, less Turkish-focused corpus composition.

Overall, POS tagging performance appears saturated across both scales, with nearly all base models exceeding 94 F1 and only marginal gains from scaling. SindBERT maintains parity with top-tier baselines, confirming that syntactic coverage in Turkish is largely solved for transformer-based encoders.

Named Entity Recognition The best base-scale performance is reached by BERTurk_{32k} (94.38), confirming its robustness for token-level classification. Close behind are ConvBER-Turk (94.03) and BERTurk_{128k} (93.81), while SindBERT_{base} achieves a solid 93.19, comparable to ELECTRA_{base} (93.49) and XLM-R_{base} (92.9). This indicates that SindBERT's RoBERTa-like setup neither clearly surpasses nor lags behind the most established Turkish encoders, suggesting that the NER task may already be approaching an upper limit with current dataset size and annotation quality.

At the large scale, XLM-R_{large} slightly leads (94.44), followed closely by SindBERT_{large} (93.64). Given that XLM-R was trained on over 2 TB of multilingual text, this narrow margin underscores the efficiency of SindBERT's more compact, Turkish-focused pretraining corpus.

In general, NER results reveal minimal separation between base and large encoders, indicating that model size has limited impact once sufficient Turkish data are used. SindBERT performs on par with the strongest monolingual models, underscoring the stability of its representations across token-level semantic tasks.

Offensive Language Detection For offensive language classification (OffensEval-TR 2020), we observe more pronounced differences between architectures. ConvBERTurk reaches the highest macro-F1 among base models (81.99), with ConvBERTurk_{mC4} (81.90) and BERTurk_{128k} (81.77) performing almost identically. ELEC-TRA variants and SindBERT_{base} (81.14) cluster slightly below, while distilled and multilingual models trail more clearly. These results highlight that models trained on monolingual Turkish corpora still offer clear advantages for pragmatic and domain-sensitive tasks. SindBERT_{base} thus performs solidly but not at the very top, suggesting that further pre-training on informal or social-media text could enhance its stylistic robustness.

In the large model group, SindBERT_{large} again performs best (82.29), surpassing XLM-R_{large} (81.99) and far exceeding EuroBERT_{610M} (75.57). This consistent lead across two of four downstream tasks emphasizes SindBERT's balanced architecture and effective use of Turkish-specific corpora.

TURBLIMP Table 5 reports the detailed TUR-BLIMP results for all base and large models. Overall, SindBERT_{base} achieves an average score of 90.3, which is comparable to ELECTRA_{base} and ELECTRA_{mC4} (both 89.9), while trailing behind the strongest baselines BERTurk_{32k} (93.8) and BERTurk_{128k} (95.1). A closer look at the perphenomenon results shows that SindBERT_{base} is particularly strong on *scrambling*, *suspended affixation*, *subject agreement*, and *irregular forms* (all ≥98), which are central morphosyntactic phenomena of Turkish. At the same time, it struggles with *ellipsis* (59.0) and *island effects* (64.0), two categories that remain challenging across most models.

For the large models, SindBERT_{large} reaches an average of 89.8, placing it slightly below EuroBERT_{610M} (90.0) and XLM-R_{large} (92.7). Its strengths mirror the base variant: ceiling-level performance in morphologically rich categories such as *suspended affixation*, *scrambling*, and *irregular forms*. However, SindBERT_{large} shows a severe weakness in *ellipsis* (27.8), which strongly lowers its overall average.

These findings highlight that monolingual models like SindBERT capture Turkish-specific morphosyntax particularly well, while multilingual models such as XLM-R generalize more effectively to harder syntactic phenomena (e.g., ellipsis and binding). This suggests a trade-off between specialization in language-specific structures and broader generalization capacities learned from multilingual corpora.

5 Discussion

5.1 Principal Findings

Our evaluation shows that SindBERT_{base} performs competitively with other widely used Turkish encoders, confirming the robustness of its RoBERTa-style pretraining setup. At the same time, SindBERT_{large} achieves the best overall results in two of four downstream tasks, notably in part-of-speech tagging and offensive language detection, and also performs strongly on several linguistic control tests. While scaling does not produce uniform gains across all benchmarks, these task-specific improvements suggest that larger contextual capacity primarily benefits pragmatically and syntactically complex settings. Similar saturation effects are visible for EuroBERT and XLM-R, indicating that many Turkish benchmarks may no longer be sufficiently discriminative to reveal consistent scaling trends. Nonetheless, diagnostic evaluations such as TURBLIMP underscore Sind-BERT's strengths in Turkish-specific grammatical phenomena (e.g., scrambling, suspended affixation, subject agreement), highlighting the model's linguistic depth beyond aggregate scores.

5.2 Corpora

A likely factor explaining the limited scaling gains lies in the training corpus composition. SindBERT was trained on 312 GB of text—dominated by mC4 (242 GB), which provides broad coverage but is considerably noisier than smaller, curated datasets. By contrast, BERTurk, trained on only a fraction of that volume but sourced from cleaner collections (OSCAR, Wikipedia, OPUS, and non-public), achieves excellent results, particularly on linguistically sensitive evaluations. This mirrors trends observed in other monolingual models such as Gott-BERT, CamemBERT, and GeistBERT, where performance gains stemmed not merely from data size but from an effective balance of quality, domain diversity, and linguistic representativeness. Our findings therefore reinforce that corpus curation, not scale alone, is decisive for progress in Turkish NLP.

Model	PoS	WikiANN	OffensEval-TR 2020	TURBLIMP AVG
ELECTRA _{small}	94.28	91.92	78.17	80.6
DistilBERTurk	94.01	91.54	79.19	87.2
ConvBERTurk	94.41	94.03	81.99	60.8
ConvBERTurk _{mC4}	94.57	93.56	<u>81.90</u>	55.5
ELECTRA _{base}	94.29	93.49	81.54	89.9
ELECTRA _{mC4}	94.4	93.43	81.38	89.9
BERTurk _{32k}	93.16	94.38	81.03	<u>93.8</u>
RoBERTurk	87.99	81.09	70.01	-
SindBERT _{base}	<u>94.47</u>	93.19	81.14	90.3
$mmBERT_{small}$	93.75	92.51	77.28	85.1
BERTurk _{128k}	94.44	93.81	81.77	95.1
EuroBERT _{210M}	92.97	90.91	75.73	86.3
XLM-R _{base}	94.23	92.9	79.77	89.2
$mmBERT_{base}$	93.75	93.35	78.49	89.3
SindBERT _{large}	94.63	93.64	82.29	89.8
XLM-R _{large}	94.39	94.44	<u>81.99</u>	92.7
EuroBERT _{610M}	93.33	91.85	75.57	90.0

Table 4: Evaluation results across four Turkish downstream tasks. Best results are shown in bold and second-best results are underlined, with rankings reported separately for base and large model groups. For the 13 base models, third-best results are additionally marked with a dotted underline. **PoS**: micro-F1 on concatenated UD datasets. **NER**: entity-level F1 on WikiANN Turkish. **Sentiment**: macro-F1 on OffensEval-TR 2020. **TurBLiMP**: average accuracy over 16 linguistic acceptability phenomena. Reported scores for PoS, NER and classification are computed on the test set, with the best checkpoint per model—task combination selected based on validation performance. TurblimP was evaluated using its predefined configuration.

Model	Ana. Agr.	Arg. Tr.	Arg. Ditr.	Bind.	Det.	Ellip.	Irr.	Isl.	Nom.	NPI	Pass.	Quant.	RelCl.	Scramb.	Subj. Agr.	Susp. Aff.	AVG
ELECTRA _{small}	74.1	86.6	79.3	70.7	91.8	10.6	98.7	39.1	90.0	90.9	100.0	97.9	79.9	99.5	82.8	97.5	80.6
DistilBERTurk	<u>96.9</u>	97.5	95.4	93.0	82.9	13.6	94.1	47.4	95.6	92.1	98.8	<u>98.4</u>	92.0	99.8	97.0	100.0	87.2
ConvBERTurk	34.3	41.9	68.1	87.4	0.0	40.5	91.2	99.3	55.6	81.5	100.0	99.0	50.9	55.9	35.7	30.9	60.8
ConvBERTurk _{mC4}	40.7	49.9	43.2	0.3	0.0	34.7	84.1	<u>95.5</u>	67.3	88.1	100.0	99.0	49.1	47.9	46.3	41.5	55.5
ELECTRA _{base}	94.3	99.6	96.1	96.2	99.3	49.7	97.9	35.3	96.6	96.1	91.2	98.0	90.7	100.0	99.0	99.0	89.9
ELECTRA _{mC4}	94.3	99.4	95.5	91.4	98.2	46.3	99.0	41.8	97.0	95.0	93.6	98.0	92.0	100.0	97.2	99.1	89.9
BERTurk _{32k}	96.7	99.7	99.8	99.9	99.9	87.4	98.8	49.4	97.4	98.2	82.2	95.7	97.7	100.0	98.3	100.0	93.8
RoBERTurk	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
SindBERT _{base}	93.7	98.3	92.9	94.6	94.2	59.0	98.0	64.0	93.9	88.7	84.8	98.3	89.9	100.0	94.0	100.0	90.3
$mmBERT_{small}$	73.3	87.6	86.5	64.2	92.9	65.8	91.3	65.1	90.2	81.0	90.1	93.9	88.5	99.4	92.3	99.0	85.1
BERTurk _{128k}	97.3	99.8	96.1	97.7	99.0	96.6	99.9	60.9	98.9	97.2	84.4	98.0	97.0	99.6	99.1	100.0	95.1
EuroBERT _{210M}	88.3	86.5	83.2	82.7	98.8	48.6	98.5	55.7	89.3	84.8	90.0	94.4	90.0	100.0	91.9	98.0	86.3
XLM-R _{base}	94.6	91.4	89.2	92.7	98.6	65.9	92.8	52.4	91.4	92.4	90.4	96.0	84.6	100.0	95.0	<u>99.7</u>	89.2
$mmBERT_{base}$	85.2	91.5	93.6	86.2	94.7	82.5	96.3	58.9	92.5	84.8	93.1	93.4	83.6	99.6	94.2	99.4	89.3
SindBERT _{large}	90.4	98.4	91.6	95.3	97.1	27.8	100.0	67.2	95.3	91.1	94.8	98.8	94.4	100.0	94.9	100.0	89.8
XLM-R _{large}	94.7	96.5	96.7	98.5	98.7	86.9	94.8	68.8	91.6	91.3	85.0	93.4	91.4	99.8	95.8	100.0	92.7
EuroBERT _{610M}	90.1	96.6	92.4	92.1	95.8	78.3	95.9	53.0	94.1	84.9	87.9	92.2	92.0	100.0	<u>95.7</u>	<u>99.5</u>	90.0

Table 5: Detailed TURBLIMP evaluation across 16 linguistic acceptability phenomena. Best results are shown in bold and second-best results are underlined, with rankings reported separately for base and large model groups. For the 13 base models, third-best results are additionally marked with a dotted underline.

A further dimension concerns vocabulary design. SindBERT employs a 52k BPE vocabulary that balances coverage and efficiency, whereas BERTurk also released a 128k-token variant, which ranks among the strongest performers in our benchmarks, especially on TURBLIMP. Recent work by Toraman et al. (2023) corroborates that vocabulary size has a substantial impact on Turkish models due to the language's agglutinative morphology. They report that optimal vocabulary scales differ by tokenization strategy: for BPE or WordPiece, vocabularies around 20% of model parameters tend to be most effective, while morphological or word-level tokenizers may benefit from substantially larger ratios. Our results align with this observation: BERTurk_{128k} profits from an expanded vocabulary despite its smaller corpus, whereas SindBERT's 52k vocabulary remains sufficiently expressive to achieve competitive results given its broader but noisier training data.

5.3 Efficiency

From an efficiency perspective, our findings highlight a favorable trade-off between scale and performance. While SindBERT_{base} achieves results comparable to its larger counterpart at a fraction of the computational cost, SindBERT_{large} still demonstrates measurable advantages on more demanding or pragmatically complex tasks. This indicates that the large model's additional capacity is not wasted, but rather contributes selectively where richer contextual representations are required. Nevertheless, for most real-world scenarios, the base configuration offers an excellent balance between efficiency and accuracy. Taken together, the flat scaling behavior across multiple Turkish model families suggests that future progress will hinge less on parameter growth and more on corpus quality, tokenization, and task design.

6 Future Directions

Future work may extend SindBERT in several directions. First, while GeistBERT built on the Gott-BERT checkpoint through continued pre-training on in-domain data (Scheible-Schmitt and Frei, 2025), and ChristBERT explored the effects of continued pre-training versus training from scratch using both general and domain-specific vocabularies, a similar ablation study has not yet been conducted for Turkish. SindBERT provides a natural starting point for replicating these approaches,

enabling systematic comparisons of domain adaptation strategies in Turkish.

Second, our findings indicate that many existing benchmarks are already saturated, as they fail to reveal consistent improvements from larger models. To overcome this limitation, future evaluations should adopt more comprehensive and discriminative test suites. In particular, the recently released TrGLUE benchmark¹ offers a promising step in this direction, providing a diverse collection of tasks. It includes natural language inference, paraphrase detection, sentiment analysis, and question answering, that more closely mirror the breadth of the original GLUE suite. Incorporating TrGLUE into future experiments would enable a more finegrained assessment of SindBERT's generalization capabilities across both syntactic and semantic dimensions.

Third, extending evaluation to specialized domains such as biomedical or legal language remains an important frontier for Turkish NLP, where Sind-BERT could serve as a foundation for targeted domain adaptation, just as GottBERT (Scheible et al., 2024) and GeistBERT (Scheible-Schmitt and Frei, 2025) did for ChristBERT (He et al., 2025).

Finally, future pre-training efforts could further improve linguistic coverage by considering document or sentence boundaries during sampling and by employing WWM (Martin et al., 2020; Chan et al., 2020).

7 Conclusion

We introduced SindBERT, the first large-scale RoBERTa encoder trained from scratch on 312 GB of Turkish text. Across four benchmarks, it performs competitively with existing models, with SindBERT_{large} achieving the best results in two tasks. While scaling brings only selective gains, this mirrors trends in XLM-R and EuroBERT, suggesting that Turkish benchmarks are nearing saturation. The contrast with BERTurk highlights the decisive role of corpus quality and variance over size. Together, these findings show that progress in Turkish NLP will depend less on scaling and more on curated data, adaptive tokenization, and challenging evaluation suites. As the first openly released large-scale RoBERTa model for Turkish, SindBERT establishes a solid foundation for future Turkish NLP.

https://huggingface.co/datasets/
turkish-nlp-suite/TrGLUE

Limitations

This work has several limitations. First, SindBERT was trained on three large-scale Turkish corpora (mC4, OSCAR23, Wikipedia) with only light filtering applied, restricted to the removal of documents containing invalid character encodings. No additional cleaning, quality filtering, or cross-source deduplication was performed. As a result, residual noise, duplicated content, and potential biases are likely to remain in the training data and may influence the learned representations.

Second, the training data was drawn exclusively from web-based sources, without explicit control for dialectal or register variation (e.g., Ottoman vs. Modern Turkish, formal vs. colloquial, or regional varieties). This may limit the model's robustness on underrepresented varieties or in specialized domains such as biomedical or legal text, unless additional domain-adaptive pre-training is performed.

Third, SindBERT was pre-trained with conservative hyperparameter settings and without extensive exploration of alternative masking strategies (e.g., Whole Word Masking) or longer training schedules. Pre-training was also conducted without mixed precision, which increased computational cost and limited the feasibility of scaling to larger model sizes or more training steps.

Fourth, we did not perform a systematic error analysis of downstream results. Such an analysis could provide insights into systematic weaknesses (e.g., frequent PoS confusions, NER boundary errors, sentiment misclassifications, or TURBLIMP minimal pair failures) and help prioritize future improvements in model design and dataset composition.

Fifth, baseline reproducibility introduces some uncertainty. ConvBERTurk and ConvBERTurk_{mC4} are based on the ELECTRA codebase, but during conversion from the original checkpoints to HuggingFace Transformers the distinction between generator and discriminator is not explicit. While ELECTRA's conversion script allows specifying this choice, ConvBERTurk appears to default to the discriminator. This may not invalidate comparisons, but it does leave open the possibility of subtle architectural differences and explains the suboptimal performance on TurblimP.

Lastly, our evaluation focused on four downstream tasks (PoS tagging, NER, sentiment classification, TURBLIMP). While these cover a diverse range of morphosyntactic, semantic, and syntactic phenomena, they do not capture the full scope of Turkish NLP challenges such as question answering, natural language inference, summarization, or long-context understanding. The generalization of SindBERT to these settings remains to be established.

Ethical Considerations

Like all large-scale language models, SindBERT may inherit biases from its training data, which can influence downstream tasks such as classification or decision-making. While no deduplication was applied, the corpus may still contain redundancy and noise, as well as deeper societal or representational biases. Furthermore, training on large webbased corpora raises privacy concerns, as models may inadvertently retain sensitive information. Responsible deployment is especially important in high-stakes domains like legal, medical, or financial NLP.

Despite optimizations for efficiency, pre-training and evaluating transformer models remain computationally demanding, contributing to energy use and carbon emissions. These environmental costs highlight the need for balancing model performance with sustainable development goals.

Acknowledgments

The authors gratefully acknowledges the support of Google's TPU Research Cloud for providing access to Cloud TPUs, which enabled efficient pretraining of SindBERT. The authors also thank Nora Limbourg, the assigned Google Cloud Customer Engineer, for her valuable technical assistance and coordination throughout the project. Finally, the authors gratefully acknowledge the scientific support and resources of the AI service infrastructure LRZ AI Systems provided by the Leibniz Supercomputing Centre (LRZ) of the Bavarian Academy of Sciences and Humanities (BAdW), funded by Bayerisches Staatsministerium für Wissenschaft und Kunst (StMWK).

References

Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. 2019. FLAIR: An Easy-to-Use Framework for State-of-the-Art NLP. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations), pages 54–59, Minneapolis, Minnesota. Association for Computational Linguistics.

- Ezgi Başar, Francesca Padovani, Jaap Jumelet, and Arianna Bisazza. 2025. Turblimp: A turkish benchmark of linguistic minimal pairs. Preprint, arXiv:2506.13487.
- Nicolas Boizard, Hippolyte Gisserot-Boukhlef, Duarte M. Alves, André Martins, Ayoub Hammal, Caio Corro, Céline Hudelot, Emmanuel Malherbe, Etienne Malaboeuf, Fanny Jourdan, Gabriel Hautreux, João Alves, Kevin El-Haddad, Manuel Faysse, Maxime Peyrard, Nuno M. Guerreiro, Patrick Fernandes, Ricardo Rei, and Pierre Colombo. 2025. Eurobert: Scaling multilingual encoders for european languages. Preprint, arXiv:2503.05500.
- Çağrı Çöltekin. 2020. A corpus of turkish offensive language on social media. In <u>Proceedings</u> of The 12th Language Resources and <u>Evaluation</u> Conference, pages 6174–6184, Marseille, France.
- Branden Chan. 2020. XLM-RoBERTa: The multilingual alternative for non-english NLP. Library Catalog: towardsdatascience.com.
- Branden Chan, Stefan Schweter, and Timo Möller. 2020.
 German's next language model. In Proceedings of the 28th International Conference on Computational Linguistics, pages 6788–6796, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. Preprint, arXiv:2003.10555.
- Pieter Delobelle, Thomas Winters, and Bettina Berendt. 2020a. RobBERT: a Dutch RoBERTa-based Language Model. arXiv:2001.06286 [cs]. ArXiv: 2001.06286.
- Pieter Delobelle, Thomas Winters, and Bettina Berendt. 2020b. RobBERT: a Dutch RoBERTa-based Language Model. In Findings of the Association for Computational Linguistics: EMNLP 2020, pages 3255–3265, Online. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. 2020. Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. Preprint, arXiv:2002.06305.

- Henry He, Johann Frei, and Raphael Scheible-Schmitt. 2025. The word and the way: Strategies for domain-specific BERT pre-training in german medical NLP. ISSN: 2693-5015.
- Tim Jansen, Yangling Tong, Victoria Zevallos, and Pedro Ortiz Suarez. 2022. Perplexed by quality: A perplexity-based method for adult and harmful content detection in multilingual heterogeneous web data. Preprint, arXiv:2212.10440.
- Zihang Jiang, Weihao Yu, Daquan Zhou, Yunpeng Chen, Jiashi Feng, and Shuicheng Yan. 2021. Convbert: Improving bert with span-based dynamic convolution. Preprint, arXiv:2008.02496.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, XiaoChen, Linlin Li, Fang Wang, and Qun Liu. 2020.Tinybert: Distilling bert for natural language understanding. Preprint, arXiv:1909.10351.
- Norman P. Jouppi, George Kurian, Sheng Li, Peter Ma, Rahul Nagarajan, Lifeng Nai, Nishant Patil, Suvinay Subramanian, Andy Swing, Brian Towles, Cliff Young, Xiang Zhou, Zongwei Zhou, and David Patterson. 2023. TPU v4: An optically reconfigurable supercomputer for machine learning with hardware support for embeddings. Preprint, arxiv:2304.01433 [cs].
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open Source Toolkit for Statistical Machine Translation. In Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv:1907.11692 [cs]. ArXiv: 1907.11692.
- Marc Marone, Orion Weller, William Fleshman, Eugene Yang, Dawn Lawrie, and Benjamin Van Durme. 2025. mmbert: A modern multilingual encoder with annealed language learning. Preprint, arXiv:2509.06888.
- Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. 2020. CamemBERT: a Tasty French Language Model. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7203–7219, Online. Association for Computational Linguistics.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael

Auli. 2019. fairseq: A Fast, Extensible Toolkit for Sequence Modeling. arXiv:1904.01038 [cs]. ArXiv: 1904.01038.

Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Crosslingual name tagging and linking for 282 languages. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1946–1958, Vancouver, Canada. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. OpenAI Blog, 1(8):9.

Afshin Rahimi, Yuan Li, and Trevor Cohn. 2019. Massively multilingual transfer for NER. In <u>Proceedings</u> of the 57th Annual Meeting of the <u>Association</u> for Computational Linguistics, pages 151–164, Florence, Italy. Association for Computational Linguistics.

Raphael Scheible, Johann Frei, Fabian Thomczyk, Henry He, Patric Tippmann, Jochen Knaus, Victor Jaravine, Frank Kramer, and Martin Boeker. 2024. GottBERT: a pure German language model. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 21237–21250, Miami, Florida, USA. Association for Computational Linguistics.

Raphael Scheible-Schmitt and Johann Frei. 2025. Geistbert: Breathing life into german nlp. <u>Preprint</u>, arXiv:2506.11903.

Stefan Schweter. 2020. Berturk - bert models for turkish. https://doi.org/10.5281/zenodo. 3770924. Version 1.0.0, Zenodo.

Stefan Schweter. 2025. BERTurk v2. https://doi.org/10.5281/zenodo.14963493. Version 2.0.0, Zenodo.

Nuri Tas. 2024. Roberturk: Adjusting roberta for turkish. Preprint, arXiv:2401.03515.

Cagri Toraman, Eyup Halit Yilmaz, Furkan Şahinu,c, and Oguzhan Ozcelik. 2023. Impact of tokenization on language models: An analysis for turkish. <u>ACM Transactions on Asian and Low-Resource Language Information Processing</u>, 22(4):1–21.

Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. BLiMP: The benchmark of linguistic minimal pairs for English. Transactions of the Association for Computational Linguistics, 8:377–392.

A Runtime

Table 7 lists the hyperparameters of the best Sind-BERT models (selected by validation performance)

for each benchmark, supporting reproducibility of our results. For transparency, Table 6 reports the total computation time per task, showing that all Turkish downstream experiments together required roughly 425 GPU hours (about 17.7 days). All base model experiments were run on an NVIDIA RTX 3090, and large model experiments on an NVIDIA H100 GPU.

TURBLIMP is not reported, as the pipeline did not record training time. Since no hyperparameter search was involved, this omission is minor and corresponds to only a few additional hours.

Task	Computation Time
PoS	200:21
WikiANN	131:02
OffensEval-TR 2020	93:37
Total	425:01

Table 6: Computation time in hours and minutes for the Turkish downstream tasks, summing to about 425 hours and 1 minute (approximately 17.7 days).

Model	Pos	S	NE	R	Sentiment		
1/10401	BF	LR	BF	LR	BF	LR	
ELECTRA _{small}	5e-05	32	5e-05	16	2e-05	16	
DistilBERTurk	2e-05	16	5e-05	16	7e-06	16	
ConvBERTurk	5e-05	32	1e-05	16	7e-06	16	
ConvBERTurk _{mC4}	5e-05	32	2e-05	32	5e-06	32	
ELECTRA _{base}	5e-05	16	2e-05	32	7e-06	32	
BERTurk _{32k}	2e-05	32	2e-05	16	7e-06	16	
RoBERTurk	5e-05	16	2e-05	16	1e-05	32	
SindBERT _{base}	1e-05	16	1e-05	32	2e-05	32	
$mmBERT_{small}$	5e-05	32	2e-05	16	2e-05	32	
BERTurk _{128k}	7e-06	16	5e-05	32	7e-06	32	
EuroBERT _{210M}	7e-06	16	5e-06	16	1e-05	32	
XLM-R _{base}	5e-06	16	1e-05	16	7e-06	16	
SindBERT _{large}	1e-05	32	7e-06	16	7e-06	16	
XLM-R _{large}	7e-06	16	5e-06	16	7e-06	32	
EuroBERT _{610M}	1e-05	16	5e-06	32	5e-06	32	

Table 7: Hyperparameters of the best-performing downstream task model for each pre-trained model. $\bf BF$ denotes the batch size, $\bf LR$ the learning rate.