

Phonikud: Hebrew Grapheme-to-Phoneme Conversion for Real-Time Text-to-Speech

Yakov Kolani¹

Maxim Melichov²

Cobi Calev¹

Morris Alper³

¹Independent Researcher

²Reichman University

³Tel Aviv University

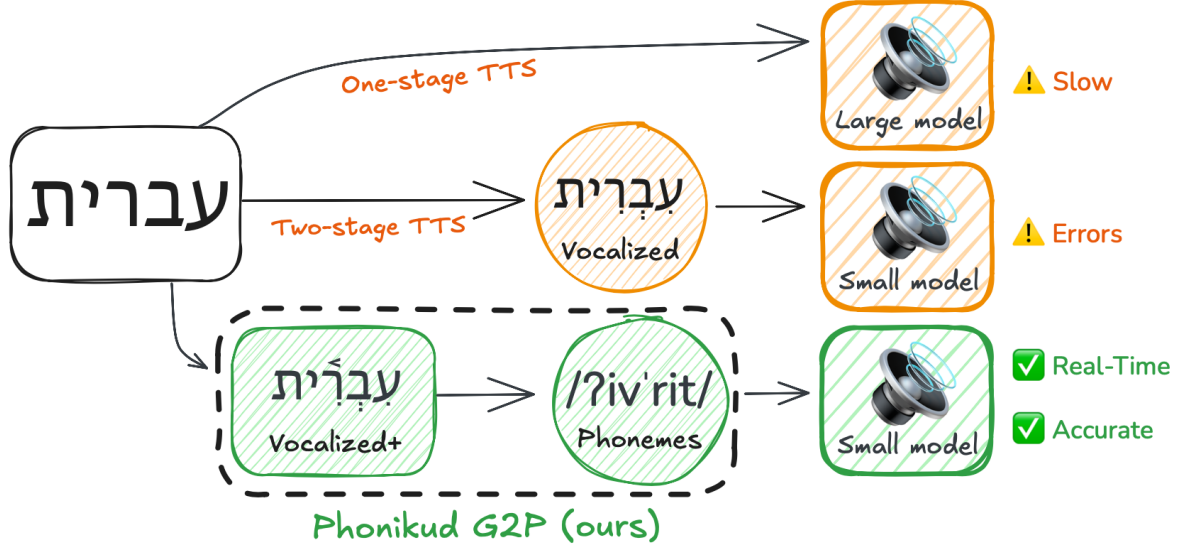


Figure 1: **Phonikud: Hebrew G2P conversion for fast, phonetically accurate Hebrew TTS.** Hebrew writing normally omits vowels, creating a speed-accuracy tradeoff for TTS: large models (orange, top) trained on raw Hebrew text are slow, while small models trained on vocalized text (orange, middle) produce pronunciation errors since added vowel marks still omit critical features like stress. Phonikud (green) achieves fast and phonetically accurate TTS via enhanced vocalization (*Vocalized+* above), augmenting standard vowel marks with additional symbols indicating phonetic features such as stress, followed by conversion to standard IPA phonemes. This enables training small TTS models with phonetically accurate outputs suitable for real-time applications.

Abstract

Real-time text-to-speech (TTS) for Modern Hebrew is challenging due to the language’s orthographic complexity. Existing solutions ignore crucial phonetic features such as stress that remain underspecified even when vowel marks are added. To address these limitations, we introduce *Phonikud*, a lightweight, open-source Hebrew grapheme-to-phoneme (G2P) system that outputs fully-specified IPA transcriptions. Our approach adapts an existing diacritization model with lightweight adaptors, incurring negligible additional latency. We also contribute the *ILSpeech* dataset of transcribed Hebrew speech with IPA annotations, serving as a benchmark for Hebrew G2P and as training data for TTS systems. Our results demonstrate that Phonikud G2P conversion more accurately predicts phonemes from Hebrew text compared to prior methods, and that this enables training of effective real-time Hebrew TTS models

with superior speed-accuracy trade-offs. We release our code, data, and models at <https://phonikud.github.io>.

1 Introduction

Despite the Modern Hebrew language being spoken by approximately nine million people (Lewis, 2009), it currently lacks an open-source real-time text-to-speech (TTS) system with adequate performance. TTS systems for important applications such as screen readers for visually impaired users and for smart home technology must run locally in real-time on resource-constrained devices. However, applying standard techniques to Hebrew is challenging due to the language’s opaque orthography, which is difficult to parse directly for the small TTS models needed to achieve low latency.

The Hebrew script omits phonetic features such as vowel sounds, leaving them to be inferred from

context. For instance, in Hebrew the word ספר may be read as /sefer/ (“book”), /sa¹par/ (“barber”), /sa¹far/ (“he counted”), or /sfar/ (“suburb”). A system of optional diacritics (*nikud*) may be used to indicate these features, but they are mostly confined to pedagogical texts such as dictionaries. Moreover, the pronunciation of a Hebrew word cannot be unambiguously determined even when vowel diacritics are provided. For example, בִּירָה may be read as either /bira/ (“beer”) or /bi¹ra/ (“capital city”). This *phonetic underspecification* challenges TTS systems, which must receive normal (unvocalized) Hebrew text as input and output correctly-pronounced Hebrew audio.

One approach maps unvocalized Hebrew text directly to audio (Roth et al., 2024; Zeldes et al., 2025). However, the large models needed to capture the complexities of Hebrew orthography incur high latency, making them unsuitable for real-time applications. Conversely, small TTS models struggle to predict accurate pronunciation from unvocalized Hebrew. Existing approaches predict vowel diacritics directly (Sharoni et al., 2023; Pratap et al., 2024), but this does not fully resolve ambiguity (as in the example above), leading to inaccurate pronunciations in TTS outputs.

To bridge this gap, we propose a lightweight grapheme-to-phoneme (G2P) pipeline, *Phonikud*, to resolve the phonetic ambiguities in written Hebrew. Specifically, we adapt an existing state-of-the-art (SOTA) model for predicting Hebrew vowel diacritics (Shmidman et al., 2023), adding lightweight adaptors to efficiently predict additional phonetic features such as stress and *shva* realization (see Section 2) needed for disambiguation. A rule-based module converts these outputs into the International Phonetic Alphabet (IPA). We show that this allows effectively training small, real-time-capable TTS models. During inference, these models accept IPA input directly for precise phonetic control; additionally, they may be applied to unvocalized Hebrew text by applying our G2P conversion with *Phonikud*.

As an additional step to this goal, we contribute the novel *ILSpeech* dataset and benchmark, consisting of high-quality Hebrew speech recordings along with Hebrew text and expert-annotated IPA transcriptions. This serves both as an additional training resource for Hebrew TTS, which currently has a dire lack of available open data, as well as a benchmark for evaluating the novel task of G2P for Hebrew text.

In summary, our key contributions are:

- A lightweight, open-source G2P model augmenting an existing Hebrew diacritizer to accurately transcribe Hebrew text in IPA.
- Results demonstrating that G2P is beneficial for training real-time TTS systems for Hebrew, along with comparisons to existing systems.
- *ILSpeech*, a novel dataset and benchmark of Hebrew speech recordings, Hebrew and IPA transcriptions, enabling TTS training and benchmarking Hebrew G2P.

We release¹ our data, code, and trained models to spur development of open-source real-time Hebrew TTS systems.

2 Phonetic Underspecification in Hebrew

Hebrew is normally written without vowel marks (*unvocalized text*), but even when these are added (*vocalized text*) it is still underspecified for various phonetic features that are needed for accurate TTS. These may be split into three primary issues:

Stress. Lexical stress is only partially predictable from word shape and part of speech in Hebrew (Graf and Ussishkin, 2003). As illustrated by the minimal pair /txina/ (“tahini”) vs. /txi¹na/ (“grinding”), both spelled טָחִינָה, stress is not indicated in the orthography even when vowel marks are provided.

Shva. The vowel mark known as *shva* is polyvalent, being either silent or pronounced as /e/. Its pronunciation depends on complex morpho-phonological rules with many irregularities (Weinberg, 1966). For example, in בִּלְלוֹנְדוֹן /be¹london/ (“in London”) the shva vowel between the first two consonants is pronounced, while in בְּלוֹנְדִּינִי /blon¹dini/ (“blonde”) it is silent.

Irregular words. Infrequently, words may deviate from regular pronunciation rules. A notable case is loanwords containing the phoneme /w/, written identically to /v/. For example, פִּינְגְּוִין /pingwin/ (“penguin”) is indistinguishable from the hypothetical form */pingvin/. Other examples of irregular spellings include זָאֵלָה /jala/ (“come on”) and יִשָּׂאֲחָר /jisa¹xar/ (“Issachar”).

These ambiguities motivate our approach of augmenting existing diacritization with additional phonetic disambiguation before converting to IPA.

¹<https://phonikud.github.io>

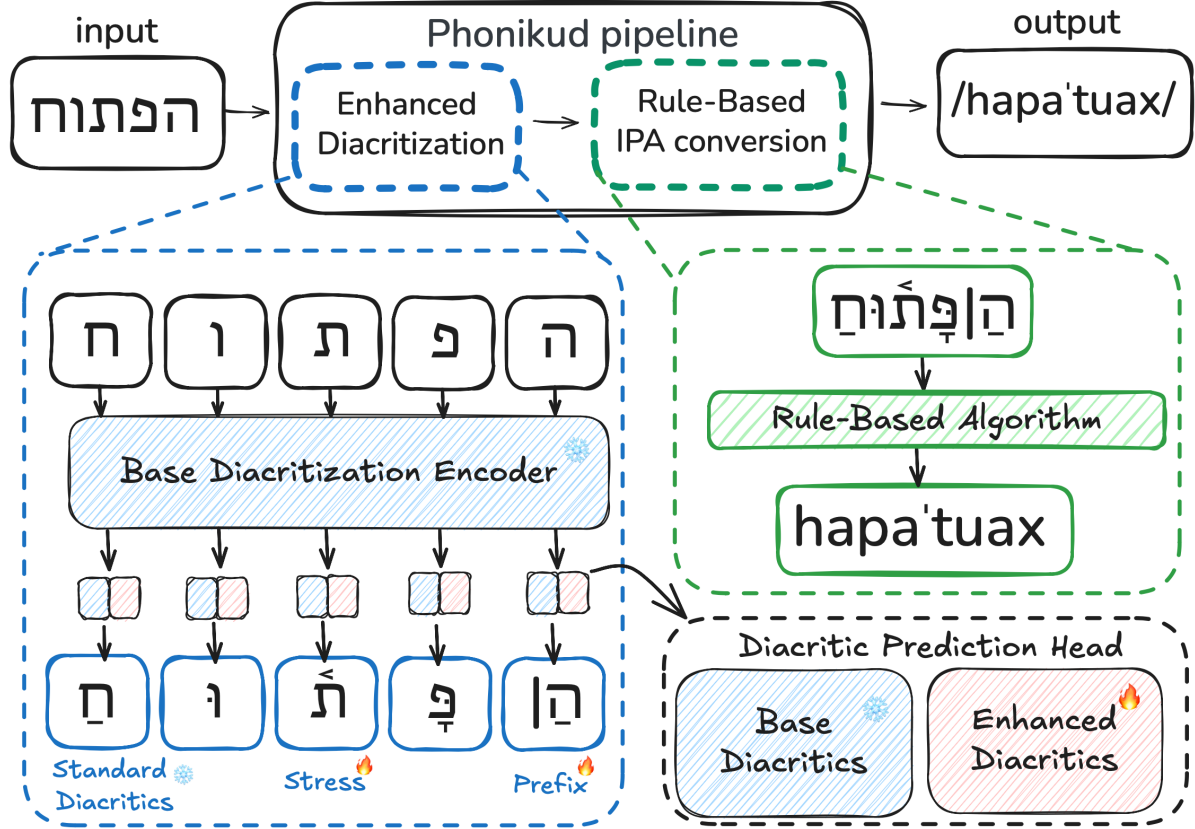


Figure 2: **The Phonikud grapheme-to-phoneme pipeline.** Phonikud converts unvocalized Hebrew text into fully-specified IPA in two steps: First, an enhanced diacritization module adds standard vowel marks and enhanced phonetic symbols to each letter of the text. This is done with a frozen (ice symbol above) base diacritization model, which is a character-level encoder model, and its per-character prediction head for standard vowel diacritics. This is augmented with an additional trainable (fire symbol above) linear adaptor serving as a head for predicting enhanced diacritics, which disambiguate the text’s phonetic content. Second, a rule-based transformation module converts this text into IPA, which may be used to train small, real-time TTS models.

3 Method

Our Phonikud system is illustrated in Figure 2. As a G2P pipeline, Phonikud takes unvocalized Hebrew text as input and outputs fully-specified IPA transcriptions, which can then be used to train efficient Hebrew TTS systems. We proceed to describe the two key components of this system – an enhanced diacritization module (Section 3.1) and rule-based IPA conversion module (Section 3.2) – followed by the novel procedure used for training the system’s learnable components (Section 3.3).

3.1 Enhanced Diacritization

Strong, efficient models have already been developed to add vowel diacritics to Hebrew text, achieving high accuracy on this standard vocalization task (Shmidman et al., 2023). Rather than learning to transcribe from scratch, we leverage this existing capability while extending it to predict ad-

ditional phonetic features needed for disambiguation. Our key insight is to augment an existing character-level encoder-based diacritization model with lightweight prediction heads to extend the set of symbols which it may predict.

We add three additional symbols that can be predicted for each character position, with logits output by new, trainable MLP prediction heads. We freeze the base encoder model and all existing prediction heads for standard vowel diacritics, while adding and training only the added weights for each new enhanced diacritic. This approach offers several advantages: training is extremely lightweight, inference predicts both standard and enhanced diacritics in parallel with minimal runtime overhead compared to the base diacritizer, and performance on standard diacritization remains constant since the base model is frozen.

We introduce three enhanced diacritics: (1) a

superscript angle indicating a stressed syllable (e.g. בֶּ֫לֶם ²), (2) a subscript line indicating a shva vowel pronounced as /e/ (e.g. מִתְּחַלֵּף ³), and (3) a vertical bar indicating the end of a cliticized prefix (e.g. קוֹדֵם ⁴). The first two directly indicate missing phonetic features, while the last aids dictionary matching of irregular words, resolving the issues from Section 2. These graphemes, which include traditional Biblical cantillation marks, are chosen because they are not used in ordinary writing.

3.2 Rule-Based IPA Conversion

After generating enhanced vocalized forms (e.g. בֶּ֫לֶם), the phonemic representation can be unambiguously determined. We apply a deterministic, rule-based algorithm to convert this to standard IPA (e.g. /lexem/). This is primarily implemented with a finite-state transducer, with states corresponding to character n-grams and output symbols corresponding to IPA phonemes. This addresses several orthographic complexities of Hebrew:

Many-to-one mappings. Multiple Hebrew graphemes frequently map to a single phoneme. For example, ט and ת both represent /t/, and three distinct vowel symbols may map to /e/.

Non-monotonic sequences. Some Hebrew words are parsed non-monotonically (not in a linear order), such as רִיחַ (“smell”) representing /'reax/ (not */'rexa/, which would be the reading in linear order).

Dual-function letters. The letters ל and ו may function as vowels or consonants depending on orthographic context. Words such as סִיּוּג /si'vug/ (“classification”) require complex logic to determine that the first ל represents /v/ while the second coalesces to the vowel /u/.

Irregular words. As discussed in Section 2, irregular words may require dictionary lookup to determine the correct pronunciation.

By addressing these complexities in the IPA conversion stage, Phonikud efficiently simplifies the representation used for TTS training. This offers several benefits: IPA provides a standardized linguistic representation that is easy to interpret and edit, and it maintains compatibility with multilingual training scenarios, while still preserving the essential phonetic information needed for TTS (shown in Section 5.3). Moreover, it allows for user choice

at inference time, as users may either input IPA directly or generate speech from Hebrew text by using the Phonikud enhanced diacritization model and/or IPA conversion.

3.3 Training Procedure

A fundamental challenge in our approach is the lack of existing ground-truth (GT) annotations for Hebrew phonetic features like lexical stress. To address this limitation, we employ a human-in-the-loop procedure to distill knowledge from existing resources along with manual refinement. In particular, we semi-automatically annotate a large-scale Hebrew corpus to indicate stress placement, prefix boundaries, and shva realization. We then distill this knowledge into our model by fine-tuning it on this pseudo-GT, which we find enables Phonikud to outperform existing baselines.

To produce large-scale data with pseudo-GT annotations, we adopt the IsraParlTweet corpus consisting of 5M lines of Hebrew text (Mor-Lan et al., 2024). We leverage Dicta’s⁵ morpho-phonological analysis API along with a set of known linguistic rules to automatically predict stress placement, prefix boundaries, and shva realization. As this procedure is frequently inaccurate, we correct many cases of errors via manual annotation, by sorting words types by frequency and correcting the most common items.

4 ILSpeech

We introduce *ILSpeech*, a high-quality Hebrew speech dataset with expert-annotated phonetic transcriptions. This dataset serves two primary purposes: (1) establishing a benchmark for Hebrew G2P systems by providing ground-truth phonetic transcriptions for systematic evaluation, and (2) supplying high-quality training data for Hebrew TTS development. Our dataset contains approximately two hours of studio-quality speech from two native Hebrew speakers (~1.5K sentences), representing a proof-of-concept that may be extended with additional speakers and content as needed.

ILSpeech provides time-aligned transcriptions with two parallel tiers: (1) unvocalized Hebrew text, and (2) expert-annotated IPA transcriptions. The latter fully specifies phonetic features such as stress that are ambiguous in vocalized Hebrew text. To the best of our knowledge, this is the first open Hebrew audio corpus containing full IPA transcrip-

²Pronounced /'lexem/ (“bread”).

³Pronounced /meti'xa/ (“stretch”).

⁴Pronounced /ha'kod/ (“the code”), with prefixed /ha-/

⁵<https://dicta.org.il>

Model	WER↓	WER ^σ ↓	CER↓	בוקר טוב
Phonikud (Ours)	0.19	0.15	0.04	'boker 'tov
Diacritizers*				
DictaBERT	0.38	0.24	0.08	bo'ker 'tov
Nakdimon	0.40	0.27	0.09	bo'ker 'tov
Multilingual G2P				
eSpeak NG	1.00	0.96	0.47	vvkr tov
Goruut	1.00	0.95	0.48	boʁɛʁ t'o:β
CharsiuG2P	1.00	0.99	0.71	bo:ʔab tē:b

Table 1: **G2P evaluation and example.** We test on ILSpeech, using unvocalized Hebrew as input and comparing to ground-truth IPA annotations. WER^σ indicates word error rate while disregarding mismatched stress. We illustrate performance on a Hebrew phrase with GT /'boker 'tov/. *Diacritizers use our IPA conversion with defaults for ambiguous features like stress.

tions, newly enabling evaluation of Hebrew G2P systems on previously unmeasurable features such as stress placement and shva realization.

The dataset addresses a critical gap in Hebrew speech resources. While several Hebrew audio corpora exist, they lack the phonetic detail necessary for G2P evaluation. In addition, existing corpora are mostly small-scale (Izre'el et al., 2001; Azogui et al., 2016; Marmorstein and Matalon, 2022; Conneau et al., 2023; Sharoni et al., 2023), while the only open large-scale corpora contain low-quality recordings without proper segmentation (Marmor et al., 2023; Turetzky et al., 2024), unsuitable for high-quality TTS training. As such, ILSpeech provides an important contribution with high-quality recording data along with expert phonetic annotations necessary for rigorous G2P evaluation.

We release ILSpeech to support open research in Hebrew speech technology, under a non-commercial license with ethical use requirements.

5 Results

Below, we provide results for G2P conversion with Phonikud and existing baselines (Section 5.1), compare downstream TTS using Phonikud to existing Hebrew TTS systems (Section 5.2), and ablate key components of our system (Section 5.3).

5.1 G2P Evaluation

Results of our Phonikud G2P system are shown in Table 1, evaluated relative to the ground-truth IPA annotations in ILSpeech. We calculate word- and character error rates (WER, CER) and WER when disregarding stress (WER^σ). We compare to two baselines: Firstly, we apply the existing SOTA

Model	WER↓	CER↓	RTF↓	# Params
Phonikud (Ours)				
Piper	0.08	0.02	0.09	20M
StyleTTS2	0.07	0.02	0.50	90M
Open Models				
MMS	0.20	0.06	0.21	36M
SASPEECH	0.11	0.04	0.16	28M
Robo-Shaul	0.08	0.04	1.58	23M
Proprietary Models				
Google	0.04	0.02	4.08	—
OpenAI	0.05	0.02	1.60	—

Table 2: **TTS Comparison.** We compare TTS models trained using our Phonikud G2P conversion (top), existing open models for Hebrew (middle), and proprietary models (bottom). Accuracy metrics (WER, CER) and latency (RTF) are calculated as described in Section 5.2. Our method yields a superior trade-off between these two dimensions.

Hebrew diacritizers DictaBERT (Shmidman et al., 2023) and Nakdimon (Gershuni and Pinter, 2022) with our IPA conversion, using reasonable defaults for ambiguous features (e.g. final stress, common in Hebrew). Secondly, we compare to existing open-source multilingual G2P libraries which ostensibly support Hebrew: eSpeak NG⁶, Goruut⁷, and CharsiuG2P⁸ (Zhu et al., 2022).

Our system outperforms all of these baselines: Our prediction of features such as stress improves performance significantly relative to existing diacritizers, while existing multilingual G2P systems are nearly unusable for Hebrew due to lack of dictionary support for common words (eSpeak NG) and extensive hallucinations in neural models (Goruut, CharsiuG2P). Errors in our method’s outputs stem from both occasional mistakes in predicting features such as stress, as well from limitations of the base diacritization model on which our model is built (see Section 7.1). In the appendix, we provide results of all of these G2P methods applied to a standard Hebrew text for visual comparison.

5.2 Downstream TTS Evaluation

We evaluate our method’s utility for downstream TTS by comparing models trained with Phonikud G2P on ILSpeech audio to baseline approaches for fast Hebrew TTS. In Table 2, we report performance of our method applied to train multiple open-source TTS architectures: a light-weight Piper⁹

⁶<https://github.com/espeak-ng/espeak-ng>

⁷<https://github.com/neurlang/goruut>

⁸<https://github.com/lingjzhu/CharsiuG2P>

⁹<https://github.com/rhasspy/piper>

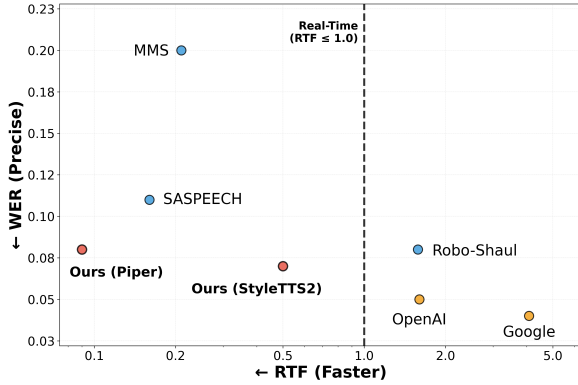


Figure 3: **Speed-accuracy trade-off.** Runtime (x-axis, log-scaled) vs. error rate (y-axis) comparison of our method (red) against open-source (blue) and proprietary (orange) TTS models. Models with $RTF \leq 1.0$ (dotted line) are real-time capable, and the lower-right direction reflects better overall performance. Our method achieves a superior speed-accuracy trade-off, with additional advantages on ambiguous features such as stress that are not reflected in automatic metrics.

Model	WER↓	CER↓
Ours	0.11	0.03
-IPA conversion	0.11	0.03
-vowel diacritics	0.24	0.09

Table 3: **Ablation study.** For ablations, we use fixed, light-weight training settings described in Section 5.3. Training on vocalized Hebrew text without IPA conversion (second row) provides comparable results, evidencing that our G2P conversion does not lose important phonetic information. However, when training on unvocalized text (last row), our small model struggles to infer correct vowels and cannot reliably produce intelligible text.

(VITS) (Kim et al., 2021) model, as well as a larger StyleTTS2 (Li et al., 2023) model. We fine-tune existing pretrained English TTS checkpoints on audio from ILSpeech along with IPA transcriptions calculated with Phonikud¹⁰. During inference, we use Phonikud to convert input Hebrew text to IPA used for synthesis.

We compare to the open-source models MMS (Pratap et al., 2024), the SASPEECH baseline (Sharoni et al., 2023), and Robo-Shaul¹¹, all light-weight models using a two-stage approach (diacritizing text followed by speech synthesis). We also compare to the proprietary TTS models

offered by Google¹² and OpenAI¹³, which support Hebrew. In the appendix, we also analyze large models which directly process unvocalized Hebrew text; these are far from real-time performance.

Following standard practice (Roth et al., 2024), we calculate error rates (WER, CER) by applying automatic speech recognition (ASR) to generations and comparing to the original input text. We also report real-time factor (RTF) values to measure latency of each system (including diacritization time, when relevant); for all open systems this is calculated on a consistent, CPU-only hardware setup to match edge computing use cases, while for proprietary models this uses their cloud inference APIs. Following prior work (Roth et al., 2024; Zeldes et al., 2025), we evaluate on a random subset of the SASPEECH (Sharoni et al., 2023) dataset, which is out-of-distribution for our model.

Our system achieves a superior trade-off between speed and accuracy than prior methods, seen visually in Figure 3. This holds both when comparing to open models run on local hardware as well as when comparing to proprietary systems run via external API, further supporting the value of our method for real-time use cases. Importantly, these automatic metrics do not capture the effect of phonetic inaccuracies such as stress placement (as ASR may still return the original, unvocalized text). Qualitative audio comparisons at our project page demonstrate our method’s improved handling of these phonetic features relative to prior methods. Our demo also illustrates that our system enables optional user control over features like stress placement.

5.3 Ablations

We ablate key parts of our system in Table 3, fixing the base TTS model (Piper) and training settings (using less training time than our main results for a light-weight comparison; see appendix for details). Removing IPA conversion (i.e. training directly on vocalized Hebrew text) yields similar objective performance to our full system. This illustrates that our IPA conversion, which has various practical advantages (Section 3.2), preserves the essential phonetic content needed to synthesize speech, while also enabling user control at inference time via either IPA or Hebrew text (with IPA conversion applied). However, training directly on undiacritized Hebrew text leads to severely degraded performance. Qualitatively, the model trained without vowels struggles

¹⁰We use IPA output from Phonikud rather than the existing manual IPA annotations to fairly evaluate our G2P pipeline and simulate scalable training on other Hebrew datasets.

¹¹<https://github.com/maxmelichov/Text-To-speech>

¹²Gemini 2.5 Flash TTS

¹³GPT-4o mini TTS

to infer the correct vowel sounds for uncommon words and often produces unintelligible output.

6 Related Work

Hebrew TTS. Early Hebrew TTS systems relied on rule-based formant synthesis (Laufer, 1975), while modern approaches are mainly learning-based. Some use two-step pipelines that add vowel diacritics before speech synthesis (Sharoni et al., 2023; Pratap et al., 2024), but this fails to resolve key phonetic ambiguities such as lexical stress placement. Others adopt an end-to-end approach, predicting speech directly from raw, undiacritized Hebrew text (Roth et al., 2024; Zeldes et al., 2025), but these require large models with high computational overhead, unsuitable for real-time use. We strike a middle ground by using a two-stage approach for computational efficiency while predicting IPA directly to ensure phonetic accuracy.

G2P conversion. Many languages have opaque orthographies, requiring TTS systems to resolve pronunciation ambiguities. Grapheme-to-phoneme conversion simplifies the learning process for TTS systems by offloading this disambiguation from the text synthesis model (Fong et al., 2019; Hexgrad, 2025). This may handle a variety of language-dependent issues, such as homograph disambiguation in English (e.g. *lead* as a verb vs. noun) (Ploujnikov and Ravanelli, 2022), predicting underspecified tone in Thai (Rugchatjaroen et al., 2019), and inferring vowels in Arabic (Elmallah et al., 2024; Kharsa et al., 2024) and Hebrew (Gershuni and Pinter, 2022; Shmidman et al., 2023). However, existing Hebrew vocalization systems and open-source G2P tools do not specify crucial phonetic features such as stress, while our method generates fully-specified phonetic transcriptions.

7 Conclusion

We have presented a new open-source Hebrew G2P system, Phonikud, and have shown that it newly enables the training of small, real-time Hebrew TTS models, which are needed for edge computing applications. Our experiments show our system compare favorably to existing solutions in both quality and runtime performance. We have also introduced the novel ILSpeech dataset and benchmark for Hebrew G2P evaluation and TTS training. We release our data, code, and trained models to enable applications and research. We envision future work building upon our contributions to further

improve Hebrew TTS performance while retaining low latency. Additional promising directions include fine-grained prosody control, support for code-switching, extensive logic for expanding symbols such as dates and addresses, semi-automated IPA annotation to increase the scale of ILSpeech, and extensions to other languages facing related phonological and orthographic challenges.

7.1 Limitations

As our method builds on existing models, we inherit a number of their limitations. The Hebrew diacritization model may output inaccurate productions, leading to incorrect IPA transcriptions. It does not support user selection among alternative vocalizations, which may be desirable in ambiguous cases. The diacritization model adheres to the conventions of formal written Hebrew, which may diverge from spoken norms (e.g. formal /sig'ri/ vs. informal /sge'ri/ for סגרי “close! (f.)”) Finally, when using our full pipeline, the prosodic quality of synthesized voice is constrained by the inherent trade-offs of real-time TTS models due to their limited capacity.

Ethics Statement

TTS is a dual-use technology: it enables valuable applications such as assistive tools for visually impaired users, but can also misused to generate disinformation or low-quality synthetic content. As with other generative models, responsible use is essential. We believe our work represents an important step towards making language technologies more accessible for lower-resourced languages such as Hebrew, while also acknowledging current limitations in representation. Our proposed dataset, like prior resources for Hebrew, cover a narrow set of speakers and styles, lacking adequate coverage of sociolinguistic variation such as the Mizrahi Hebrew accent. We anticipate that future work will increase this coverage to support more equitable and inclusive voice technologies.

Acknowledgements

We thank Dicta for encouraging our work and for approving our release of our model and data which incorporate their results. We also thank the speakers in ILSpeech for providing an essential resource for the development of Hebrew language technologies. Finally, we acknowledge Kush Jain, Shlomo Tannor, and Mark Kahn for their helpful feedback and suggestions.

References

- International Phonetic Association. 1999. *Handbook of the International Phonetic Association: A guide to the use of the International Phonetic Alphabet*. Cambridge University Press.
- Jacob Azogui, Anat Lerner, and Vered Silber-Varod. 2016. The open university of israel map task corpus (matacop).
- Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2023. Fleurs: Few-shot learning evaluation of universal representations of speech. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 798–805. IEEE.
- Muhammad Morsy Elmallah, Mahmoud Reda, Kareem Darwish, Abdelrahman El-Sheikh, Ashraf Hatim El-neima, Murtadha Aljubran, Nouf Alsaed, Reem Mohammed, and Mohamed Al-Badrashiny. 2024. Arabic diacritization using morphologically informed character-level model. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 1446–1454.
- Jason Fong, Jason Taylor, Korin Richmond, and Simon King. 2019. A comparison between letters and phones as input to sequence-to-sequence models for speech synthesis. In *The 10th ISCA Speech Synthesis Workshop*, pages 223–227. International Speech Communication Association.
- Elazar Gershuni and Yuval Pinter. 2022. Restoring hebrew diacritics without a dictionary. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1010–1018.
- Dafna Graf and Adam Ussishkin. 2003. Emergent iambs: stress in modern hebrew. *Lingua*, 113(3):239–270.
- Hexgrad. 2025. [G2p shrinks speech models](#).
- Shlomo Izre’el, Benjamin Hary, and Giora Rahav. 2001. Designing cosih: the corpus of spoken israeli hebrew. *International Journal of Corpus Linguistics*, 6(2):171–197.
- Ruba Kharsa, Ashraf Elnagar, and Sane Yagi. 2024. Bert-based arabic diacritization: A state-of-the-art approach for improving text accuracy and pronunciation. *Expert Systems with Applications*, 248:123416.
- Jaehyeon Kim, Jungil Kong, and Juhee Son. 2021. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In *International Conference on Machine Learning*, pages 5530–5540. PMLR.
- Asher Laufer. 1975. A programme for synthesizing hebrew speech. *Phonetica*, 32(4):292–299.
- M. Paul Lewis, editor. 2009. *Ethnologue: Languages of the World*, sixteenth edition. SIL International, Dallas, TX, USA.
- Yinghao Aaron Li, Cong Han, Vinay Raghavan, Gavin Mischler, and Nima Mesgarani. 2023. Stylelets 2: Towards human-level text-to-speech through style diffusion and adversarial training with large speech language models. *Advances in Neural Information Processing Systems*, 36:19594–19621.
- Yanir Marmor, Kinneret Misgav, and Yair Lifshitz. 2023. ivrit. ai: A comprehensive dataset of hebrew speech for ai research and development. *arXiv preprint arXiv:2307.08720*.
- Michal Marmorstein and Nadav Matalon. 2022. The huji corpus of spoken hebrew: An interaction-oriented design of a corpus.
- Guy Mor-Lan, Effi Levi, Tamir Sheaffer, and Shaul R Shenhav. 2024. Israparltweet: The israeli parliamentary and twitter resource. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9372–9381.
- Artem Ploujnikov and Mirco Ravanelli. 2022. Soundchoice: Grapheme-to-phoneme models with semantic disambiguation.
- Vineel Pratap, Andros Tjandra, Bowen Shi, Paden Tomasello, Arun Babu, Sayani Kundu, Ali Elkahky, Zhaoheng Ni, Apoorv Vyas, Maryam Fazel-Zarandi, and 1 others. 2024. Scaling speech technology to 1,000+ languages. *Journal of Machine Learning Research*, 25(97):1–52.
- Amit Roth, Arnon Turetzky, and Yossi Adi. 2024. A language modeling approach to diacritic-free hebrew tts. In *Proc. Interspeech 2024*, pages 2775–2779.
- Anocha Rugchatjaroen, Sittipong Saychum, Sarawoot Kongyong, Patcharika Chootrakool, Sawit Kasuriya, and Chai Wutiwiwatchai. 2019. Efficient two-stage processing for joint sequence model-based thai grapheme-to-phoneme conversion. *Speech Communication*, 106:105–111.
- Orian Sharoni, Roei Shenberg, and Erica Cooper. 2023. Saspeech: A hebrew single speaker dataset for text to speech and voice conversion. In *Proc. Interspeech*.
- Shaltiel Shmidman, Avi Shmidman, and Moshe Koppel. 2023. [Dictabert: A state-of-the-art bert suite for modern hebrew](#). *Preprint*, arXiv:2308.16687.
- Arnon Turetzky, Or Tal, Yael Segal, Yehoshua Dissen, Ella Zeldes, Amit Roth, Eyal Cohen, Yosi Shrem, Bronya R Chernyak, Olga Seleznova, and 1 others. 2024. Hebdb: a weakly supervised dataset for hebrew speech processing. In *Proc. Interspeech 2024*, pages 1360–1364.
- Werner Weinberg. 1966. Spoken israeli hebrew: Trends in the departures from classical phonology. *Journal of Semitic Studies*, 11(1):40–68.

Ella Zeldes, Or Tal, and Yossi Adi. 2025. Enhancing tts stability in hebrew using discrete semantic units. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Jian Zhu, Cong Zhang, and David Jurgens. 2022. Byt5 model for massively multilingual grapheme-to-phoneme conversion. *arXiv preprint arXiv:2204.03067*.

Appendix

A ILSpeech Dataset Details

The ILSpeech dataset consists of approximately two hours of Hebrew speech from two speakers. The speech content includes diverse topics covering science, technology, history, and everyday conversational speech. The audio was originally recorded in a studio environment at 44kHz, then enhanced using Adobe Enhance Speech v2¹⁴ and normalized to 22.05kHz. It was then segmented by automatically splitting recordings at silence boundaries with manual refinement, resulting in disjoint segments of 4-14 seconds each. Hebrew and IPA annotations were produced manually by the authors.

We release ILSpeech under a non-commercial license with ethical use requirements. For more details, please see the dataset page¹⁵.

B Implementation Details

B.1 G2P Conventions

In our enhanced vocalization scheme, the superscript angle marking stress is only used on non-final stressed symbols (since the most common stress pattern in Hebrew is final stress).

In our paper we provide IPA transcriptions using common linguistic conventions for clarity, while the implementation in our code and demo (and shown in Appendix C.3) differs slightly in the following conventions: Firstly, it uses the symbols /ɤ ʁ/ rather than /r x/ as a more narrow transcription of the sounds usually realized as uvular in Modern Hebrew. Secondly, it indicates stress immediately before the stressed vowel rather than before the entire stressed syllable. These are both illustrated in the word פִּיר, transcribed as /ɤ¹uaʁ/ or /¹ruax/ depending on convention.

¹⁴<https://podcast.adobe.com/en/enhance>

¹⁵<https://huggingface.co/datasets/thewh1teagle/ILSpeech>

B.2 G2P Pseudo-GT Construction

For pseudo-GT construction, we use the Isra-ParlTweet dataset¹⁶ of approximately 5M lines of Hebrew text from parliamentary proceedings and Tweets (Mor-Lan et al., 2024). We enhance the original text with diacritics using the Dicta system, including stress marking for words with second or third syllable stress, vertical bars to mark prefix letters, and programmatically apply vocal shva with expert editing.

B.2.1 G2P Model Architecture Details

The Phonikud architecture consists the following components: The base diacritization model is

As our base diacritization model, we use an open-weights DictaBERT model (Shmidman et al., 2023), using a checkpoint¹⁷ which has been fine-tuned for Hebrew diacritization. This is an encoder-only ~300M parameter model with linear classification heads for Hebrew diacritics.

We enhance this with an additional head for predicting additional symbols (stress, vocal shva, and prefix markers). This is implemented as a two-layer MLP with hidden dimension 256, intermediate ReLU activation, outputting logits for these new symbols. These additional layers add a negligible number of parameters to the model, allowing for efficient training and inference. They are randomly initialized and trained as described in the main paper.

B.3 G2P Training Details

We train our G2P model for approximately six epochs, until early stopping is triggered. We use a batch size of 256 and a learning rate of 5e-3, with 5% of the data reserved for validation. We train the model on a single GPU.

B.4 TTS Training Details

Our Piper-based model is initialized from a pre-trained English checkpoint¹⁸, and trained with learning rate 2e-4 and batch size 24 for ~10K epochs. Our StyleTTS2-based model is also initialized from a pre-trained English checkpoint¹⁹

¹⁶<https://huggingface.co/datasets/guymorlan/IsraParlTweet>

¹⁷https://huggingface.co/datasets/rhasspy/piper-checkpoints/resolve/main/en/en_US/ryan/medium/epoch=4641-step=3104302.ckpt

¹⁹<https://huggingface.co/dangtr0408/StyleTTS2-lite>

and trained with learning rate $1e-4$ and batch size 5 for 123 epochs. For ablations, we use a Piper model with the previous hyperparameters and training for ~ 500 epochs for a light-weight comparison. All training is conducted on a single GPU.

B.5 Evaluation Setup

To compute Word Error Rate (WER) and Character Error Rate (CER) metrics, we use a Whisper-based Hebrew ASR model²⁰, calculated only over Hebrew words. Real-Time Factor (RTF) was calculated as the ratio of processing time to audio duration (T/D), where $RTF < 1.0$ indicates real-time capability. All RTF measurements include the complete processing pipeline (including diacritization when used) for consistent comparison across models. This is conducted on standard consumer hardware (macOS M1) without GPU acceleration to reflect real-world edge computing scenarios.

We evaluate on 100 samples from SASPEECH, selected randomly from samples without special characters containing at least six words.

C Additional Results

C.1 Audio Results and Comparisons

For audio results of our model and comparisons to existing Hebrew TTS models, please refer to our project page: <https://phonikud.github.io>

C.2 Comparison to Large, One-Stage Models

The large, open-source Hebrew TTS models LoTHM (Zeldes et al., 2025) and HebTTS (Roth et al., 2024) take unvocalized Hebrew text as input directly. They exhibit significantly slower inference speeds than our method and other models we compare to (RTF of 84.75 for LoTHM and 25.44 for HebTTS, with the same hardware setup), vs. 0.09 for our method with Piper. Since $RTF > 1$ indicates performance slower than real-time, we focus our main evaluation on models that are close to meeting the performance constraints of practical edge deployment scenarios.

C.3 North Wind and Sun - Hebrew Phonetic Transcriptions

Figure 4 presents the Hebrew version of the “The North Wind and the Sun” fable reproduced from the

IPA Handbook (Association, 1999), with ground-truth phonemes, the output of our Phonikud G2P system and those of alternative G2P systems for comparison. Note that here we use the IPA conventions described in Appendix B.1.

²⁰<https://huggingface.co/ivrit-ai/whisper-large-v3-turbo-ct2>

Original Unvocalized Hebrew	רוח הצפון והשמש התוכחו ביניהם מי מהם חזק יותר. גמרו, כי את הנצחון ינחל מי שיצליח לפשוט מעל עובר אורח את בגדיו. פתח רוח הצפון ונשב בחזקה. הידק האדם את בגדיו אל גופו. אז הסתער עליו הרוח ביתר עוז, אך האדם, משהוסיף הקור לענותו, לבש מעיל עליון על בגדיו. נואש ממנו הרוח ומסרו בידי השמש. תחילה זרח עליו השמש ברכות, והאדם הסיר את בגדו העליון מעליו. הגביר השמש את חומו. עד שלא יכול האדם לעמוד בפני השרב, ופשט את בגדיו ונכנס לתוך הנהר, שהיה בקרבת מקום, כדי לרחוץ במימיו.
Ground-Truth Vocalized Hebrew	רוח הצפון והשמש התוכחו ביניהם מי מהם חזק יותר. גמרו, כי את הנצחון ינחל מי שיצליח לפשוט מעל עובר אורח את בגדיו. פתח רוח הצפון ונשב בחזקה. הידק האדם את בגדיו אל גופו. אז הסתער עליו הרוח ביתר עוז, אך האדם, משהוסיף הקור לענותו, לבש מעיל עליון על בגדיו. נואש ממנו הרוח ומסרו בידי השמש. תחילה זרח עליו השמש ברכות, והאדם הסיר את בגדו העליון מעליו. הגביר השמש את חומו. עד שלא יכול האדם לעמוד בפני השרב, ופשט את בגדיו ונכנס לתוך הנהר, שהיה בקרבת מקום, כדי לרחוץ במימיו.
Ground-Truth IPA	ʁ'uaχ hatsaf'on vehaf'emef hitvakχ'u bejneħ'em m'i meh'em χaz'ak jot'ek. gamk'u, k'i ?'et hanitsaχ'on jinχ'al m'i seǰatsl'iaχ lifʃ'ot meʔ'al ?ov'ek ?okaaχ ?'et bgad'av. pat'aχ ʁ'uaχ hatsaf'on venaf'av beχozk'a. hid'ek haʔad'am ?'et bgad'av ?'el guf'o. ?'az histaʔ'ek ?al'av haʁ'uaχ bej'etek ?'oz, ?'aχ haʔad'am, mifehos'if hak'ok laʔanot'o, lav'aʃ meʔ'il ?elj'on ?'al bgad'av. noʔ'aʃ mim'enu haʁ'uaχ umsak'o bid'ej haʃ'e-mef. tχil'a zaʁ'aχ ?al'av haʃ'emef beʁak'ut, vebaʔad'am hes'ik ?'et bigd'o haʔelj'on meʔ'al'av. highb'ik haʃ'emef ?'et χum'o, ?'ad sel'o jaχ'ol haʔad'am laʔam'od bifi'ej haʃaʁ'av, ufaʃ'at ?'et bgad'av veniχn'as let'oχ hanah'aκ, fehaj'a bekiv'at mak'om, ked'ej liʁ'χ'ots bemejm'av
Phonikud (Ours)	ʁ'uaχ hatsaf'on vehaf'emef hitukχ'u bneħ'em m'i meh'em χaz'ak jot'ek. gamk'u, k'i ?'et hanitsaχ'on jinχ'el m'i seǰatsl'iaχ lifʃ'ot meʔ'al ?ov'ek ?okaaχ ?'et bgad'av. pat'aχ ʁ'uaχ hatsaf'on venaf'av baχazk'a. hid'ek haʔad'am ?'et bgad'av ?'el guf'o. ?'az histaʔ'ek ?al'av haʁ'uaχ bjet'ek ?'oz, ?'aχ haʔad'am, mifehos'if hak'ok laʔanot'o, lav'aʃ mʔ'il ?elj'on ?'al bgad'av. noʔ'aʃ mim'enu haʁ'uaχ umask'o bid'ej haʃ'emef. tχil'a zaʁ'aχ ?al'av haʃ'emef beʁak'ot, vebaʔad'am hes'ik ?'et bgad'o haʔelj'on meʔ'al'av. highb'ik haʃ'emef ?'et χam'o, ?'ad sel'o jaχ'ol haʔad'am laʔam'od bifi'ej haʃaʁ'av, ufaʃ'at ?'et hagam'av veniχn'as let'oχ hanah'aκ, fehaj'a bekiv'at mak'om, ked'ej liʁ'χ'ots bemem'av.
Dicta*	ku'aχ hatsaf'on vefaf'emef hitukχ'u bneħ'em m'i meh'em χaz'ak jot'ek. gamk'u, k'i ?'et hanitsaχ'on jinχ'el m'i seǰatsl'iaχ lifʃ'ot meʔ'al ?ov'ek ?okaaχ ?'et bgad'av. pat'aχ ku'aχ hatsaf'on venaf'av baχazk'a. hid'ek haʔad'am ?'et bgad'av ?'el guf'o. ?'az histaʔ'ek ?al'av haʁ'uaχ bjet'ek ?'oz, ?'aχ haʔad'am, mifehos'if hak'ok laʔanot'o, lav'aʃ mʔ'il ?elj'on ?'al bgad'av. noʔ'aʃ mimen'u haʁ'uaχ umask'o bid'ej haf'emef. tχil'a zaʁ'aχ ?al'av haf'emef baχaχ'ot, vebaʔad'am hes'ik ?'et bgad'o haʔelj'on meʔ'al'av. highb'ik haf'emef ?'et χam'o, ?'ad sel'o jaχ'ol haʔad'am laʔam'od bifi'ej haʃaʁ'av, ufaʃ'at ?'et hagam'av veniχn'as lt'oχ hanah'aκ, fehaj'a bkiv'at mak'om, ked'ej liʁ'χ'ots bmem'av.
Nakdimon*	ku'aχ hatsaf'on vefaf'emef hitvokχ'u bneħ'em m'i meh'em χaz'ak jot'ek. gamk'u, k'i ?'et hanitsaχ'on jinχ'al m'i seǰatsl'iaχ lifʃ'ot meʔ'al ?ov'ek ?okaaχ ?'et bgad'av. pet'aχ ku'aχ hatsaf'on venef'ev baχaza-k'a. hid'ek haʔad'am ?'et bgad'av ?'el guf'o. ?'az histaʔ'ek ?al'av haʁ'uaχ bjet'ek ?'oz, ?'aχ haʔad'am, mifehos'if hak'ok laʔanot'o, lav'aʃ mʔ'il ?elj'on ?'al bgad'av. noʔ'aʃ mimen'u haʁ'uaχ umask'u bid'ej haf'emef. tχil'a zaʁ'aχ ?al'av haf'emef baχaχ'ot, vebaʔad'am hes'ik ?'et bgd'o haʔelj'on meʔ'al'av. highv'ik haf'am'ef ?'et χem'o, ?'ad sel'o jaχ'ol haʔad'am laʔam'od bifi'ej haʃaʁ'av, ufaʃ'at ?'et hagam'av veniχn'as lt'oχ hanah'aκ, fehaj'a bkiv'at mak'om, kd'ej liʁ'χ'ots bmejam'av.
Espeak	rvχ htsfvnə vħʃmf htvχχv vnihem mi mhem χzk joter gmrν χi ?t hntsχvnə jnχl mi ʃitsliχ lfʃvt mʔl ?vvr ?vrχ ?t vgdiv ftχ rvχ htsfvnə vnʃv vχzka hidk hʔdm ?t vgdiv ?l gvfv ?z hstʔr ?liv hrvχ vitr ?vz ?χ hʔdm mʃhvsif hkor lʔnotv lvʃ mʔil ?livnə ?l vgdiv nvʔʃ mmmν hrvχ vmsrv vjdī hʃmf tχila zrx ?liv hʃmf vrχvt vħʔdm hsr ?t vgdu hʔlivnə mʔliv hgvjr hʃmf ?t χmv ?d ʃlʔ jaχol hʔdm lʔmvd vfni hʃrv vfʃt ?t hgdiv vnχns ltvχ hnar ʃhih vkrvt mkom χdi lrxvts vmimiv
Goruut	riax hets'fun vofmaʃ hθokχau benihem mi mahem χezek jor. gemeru:, χi ?at hentsaχ'no jnχle mi ʃitsi'liχ lefʃot maʃal ʃo'ver orax et bagdi:. petax ku'aχ hets'fun onʃb βχzka. ajad'k hada:m et bagdi: al go'fu. ?az hes'tek ʃalʒov heru:x bi'ter ?oz, ?ax hada:m, mʃheosif hakur leʃnuθu, lavʃ ma'ʃil ʃalio:n ʃal bagdi:. nu'aʃ mameno heru:x omaesro bidi heʃ'mef. θi'χila zax ʃalʒov heʃ'mef braxot, vahem hesik et bagdo heʔalijon maʃali:. heg'bir heʃ'mef et xmau, ʃad ʃala jaku:l hada:m leʔumod be'fi eʃrb, opet et hegdi'o venχens to'tax hener, ʃa'ja bakarava ma'kum, χadi: lokats emi'mio.
CharsiuG2P	rôth dæʃtæ:ʌ otokomoro tata:ltʰe: bəritec mî: mætəp tʰæzõ: ze:təbo dæmbo: tøy tʰêt tariztʰol: jærtʰəl mî: kîfmitʰ lokej mâ:l tâmbæb æobʃ tʰêt bæddijo tʰatʰ rôth dæʃtæ:ʌ ɔlkʰəb betʰza: tæ:də: tæ:dəm tʰêt bæddijo ʌl dotte æz destaa:b tælio de:bo:tʰ be:təb tæz æn tæ:dən mətəde:miir dama:je:b la:re:ta lok mæ:vyl tæmie? tâl bæddijo roa:k mæmro de:bo:tʰ omməbo bidi takmakeɪ tʰatʰəda zəbth tælio takmak bæbtot: oðə:dəm tasib tʰêt bæddæ ta:mijo? mæ:v:mio tadbib: takmak tʰêt tʰəma: tæd kəla zætol tæ:dəm la:ma:d būri tatəba ofki tʰêt taddajo ɔrgələrs ləta:l tartab kedit: bo:bəbæt ma:χ tædi labtʰə? bæminio

Figure 4: Hebrew phonetic transcriptions of "The North Wind and the Sun" fable. We include our additional stress and vocal shva symbols in the vocalized Hebrew in the second row to illustrate the ground-truth pronunciation. The diacritization methods marked with an asterisk* are applied with our IPA conversion with defaults for ambiguous features like stress.