Sarika Khmer Text-to-Speech Technology

Seanghay Yath
Senior ML Engineer @ DGC
seanghay.dev@gmail.com





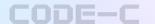


Introduction

Currently, I am a Senior ML Engineer at the Digital Government Committee. I am also a Product Developer in my free time. I built tools like KhmerScan.com and more.

My interests are product design and product development and some deep learning.



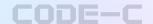






Disclaimer

The views and opinions expressed herein are solely my own and do not reflect, represent, or constitute the official position, policy, or endorsement of any organization, company, institution, or employer with which I am or have been affiliated.







What is Text to Speech

A technology that converts written text into spoken audio. It's a core component of speech synthesis systems that enables computers to "read aloud" digital text in a human-like voice.

TEXT

អ្នកស្រី កែវ ចន្ទូ បេះទុរេនលក់អស់៣០០តោនរួចហើយ
ជាមួយតម្លៃ១៥០០០រៀលក្នុងមួយគីឡូក្រាម







Why do we need it?

- **Content/News Creation**
- **Walter** Visual Impairment Support
- in Virtual Assistants / Robotics
- **Mavigation Systems**

- Video Narration
- Public Announcements
- Audiobook Production
- **Customer Service**







Human vs. TTS

Rich emotional expression and natural intonation	Consistent quality every time
Perfect contextual understanding	Available 24/7 instantly
Variable quality (fatigue, mood, health affects performance)	Can struggle with complex pronunciations
High cost and time investment	Low cost after initial setup
Difficult to modify after recording	Perfect for technical documentation and accessibility
Limited availability and scalability	Unlimited scalability







Popular TTS









IIElevenLabs







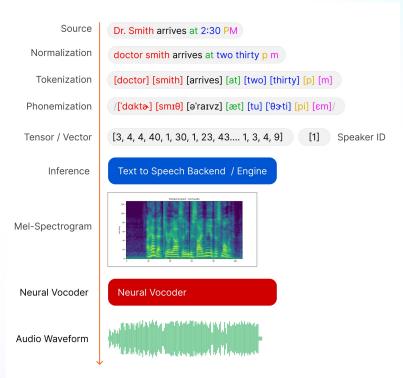


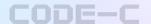






Text-to-Speech Pipeline









Khmer Language

"Khmer has been influenced considerably by Sanskrit and Pali especially in the royal and religious registers, through Hinduism and Buddhism, due to Old Khmer being the language of the historical empires of Chenla and Angkor." – Wikipedia



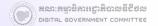




Khmer Language Challenge

- Khmer Character Encoding
- Khmer Normalization and Verbalization
- Khmer Word Boundary Tokenization
- □ Khmer Grapheme to Phoneme (Linguistic Processing)
- Khmer Pronunciation Toolkit (Linguistic Processing)







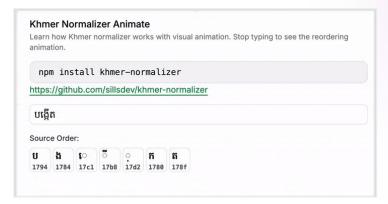
Khmer Encoding Normalization



https://ភាសាខ្មែរ.com

By Marc Durdin and Makara Sok

https://animate-khmer-normalizer.vercel.app







Discrepancies in Khmer Unicode Character Ordering Rules and a Proposed Solution

https://www.youtube.com/watch?v=mD-nrfvWtgc

13 ways to encode ស៊ើប /səəp/ 'to detect'																						
	ស ៊		៊ ើ ប			,	ស៊ើប	ų.	ស	্	េ	ី	ប	ស៊ើប		ស	េ	ិ	্	ប	ស៊ើប	
	9F	CA	ΒE	94			ııqo	⊕,		9F	вв	C1	В8	94	រេត្តប	*		C1				•
	ស	<u></u>	េ	ី	ប	,	ស៊ើប	*	ស	্	ី	េរ	ប	ស៊ើប	*	ស	ö	េ	ិ	ប	ស៊ើប	
	9F	CA	C1	В8	94				,				C1		ııqu			С9				ւսզս
*	ស		ី	េ	ប	,	ស៊ើប		•	ស	ី	্	េ	ប	ស៊ើប	*	ស	ö	ិ	េ	ប	ស៊ើប
	9F	CA	В8	C1	94		LBQO			9F	В8	вв	C1	94			9F	С9	в8	С1	94	ııqo
*	ស	্	ើ	ប		4	ស៊ើប		ស	ី	េ	্	ប	ស៊ើប	•				•			
T	9F	BB BE 94		T				BB		red O												
	ស	ើ	্	ប		1	ស៊ើប		ស	េ	্	ី	ប	ស៊ើប								
	9F	ΒE	вв	94			ııqo			9F	C1	вв	В8	94	1000							







Word Tokenization

Khmer is a language that has no clear word boundary.

For example, the sentence: "អ្នកស្រី កែវ ចន្ធ បេះទុរេនលក់អស់៣០០តោនរួចហើយ ជាមួយតម្លៃ១៥០០០រៀលក្នុងមួយគីឡូក្រាម" \rightarrow ['អ្នកស្រី', ' ', 'កែវ', ' ', 'ចន្ធ', ' ', 'បេះ', 'ទុរេន', 'លក់', 'អស់', '៣០០', 'តោន', 'រួចហើយ', ' ', 'ជាមួយ', 'តម្លៃ', '១៥០០០', 'រៀល', 'ក្នុង', 'មួយ', 'គីឡូក្រាម']

We can achieve this via khmercut Python library. Once it's broken down into piece, we can then feed it to the text processor.





Text Normalization / Verbalization

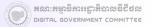
123.324 -> មួយរយ_ម្ភៃបី_ចុច_បីរយ_ម្ភៃបួន

123.001 -> មួយរយ_ម្ភៃបី_ចុច_សូន្យ_សូន្យ_មួយ

010123123 -> 0_10_12_31_23

1A 4444 -> 1 A ការ៉េ4







Text Normalization / Verbalization

(ชี) - A Khmer Text Normalization and Verbalization Toolkit

https://github.com/seanghay/tha

pip install tha

```
## Number - Cardinals
assert tha.cardinals.processor("1234") == "មួយពាន់_ពីររយ_សាមសិហ្សន"
assert tha.cardinals.processor("1") == "೪೮೮"
assert tha.cardinals.processor("1_2") == "មួយ_ពីរ"
assert tha.cardinals.processor("-1") == "ដក_មួយ"
assert tha.cardinals.processor("10") == "ដប់"
assert tha.cardinals.processor("15") == "ដប់ប្រាំ"
assert tha.cardinals.processor("100") == "មួយ ເພ"
assert tha.cardinals.processor("10000") == "មួយម៉ឺន"
assert tha.cardinals.processor("-10000.234") == "ដក មួយមីន.ពីររយ សាមសិឃន"
assert tha.cardinals.processor("-10000,234") == "ដក_មួយម៉ឺន,ពីររយ សាមសិឃន"
## Number - Decimals
assert tha.decimals.processor("123.324") == "មយរយ ម្ងៃបី ចច បីរយ ម្ដែបន"
assert tha.decimals.processor("123.001") == "មួយរយ_ម្ដែចច្រក្សន្ទ[សួន្ទ[មួយ"
assert tha.decimals.processor("-123.0012") == "ដក_មួយរយ_ម្ដែបី_ចុច_សូន្យ_សូន្យ_ដប់ពីរ"
assert tha.decimals.processor("-123,0012") == "ដក_មួយរយ_ម្ដែបី_ក្បៀស_សូន្យ_សូន្យ_ដប់ពីរ"
## Number - Ordinals
assert tha.ordinals.processor("5th") == "ອື_ຖຸຕໍ່"
assert tha.ordinals.processor("3rd") == "೯_ರ್"
assert tha.ordinals.processor("1st") == "క్ ఆటా"
assert tha.ordinals.processor("10th") == "೯ " ಟರು"
assert tha.ordinals.processor("10") == "10"
## Number - Currency
assert tha.currency.processor("$100.01") == "មួយរយដុល្លារ_មួយសេន"
assert tha.currency.processor("$100") == "មយរយ ដល្លារ"
assert tha.currency.processor("100$") == "មួយរយដុល្លារ"
assert tha.currency.processor("100៖") == "មួយរយរៀល"
assert tha.currency.processor("100.32៖") == "មួយរយ_ចូច_សាមសិបពីររៀល"
assert tha.currency.processor("100.0032៖") == "មយរយ ចច សន្យ សន្យ សាមសិបពីរវៀល"
```







Homograph Disambiguation

លោក តារា បានប្រកួតឈ្នះ លោក សុខចំនួន 4-1 តើ 4-1 ស្មើប៉ុន្មាន?

ទឹកកាន់តែឡើងបន្តិចម្តងៗ

iPhone 8, iPhone 11, iPhone 12

ស.ជ.ណ, គជប, ខ្ញុំឈឺក, លើកដៃ





Khmer Phoneme Inventory

	Bilabial	Alveolar	Palatal	Velar	Glottal
Plosives	/p/	/t/	/c/	/k/	/?/
Asp. Plosives	/ph-/	/t ^h -/	/c ^h -/	/k ^h -/	
Implosives	/b-/	/d-/			
Fricatives		/s-/			/h/
Nasals	/m/	/n/	/n/	/ŋ/	
Semi-vowels	/w/		/i/		
Lateral		/\/			
Flap		/r-/			

Khmer Character Specification/Usages - Makara Sok et. al.







Pronunciation Dictionary

111,000+ Pairs Khmer + English

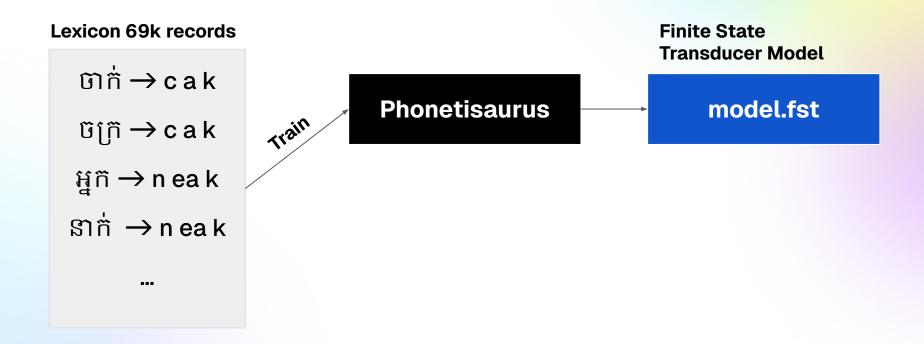
```
បង់ដា
              ban.daa
       បង់តូស
              ban.tooh
       បង់ទុយ
              ɓaŋ.tuj
        បង់បក់ baŋ.bak
        បង់បត់ ban.bat
       បង់បាញ់
              ban.ban
      បង់បោយ baŋ.baoj
        បង់ពន្
              ɓaŋ.pun
      បង់ម្សៀត
              baŋ.msiət
     បង់សូអាត
              ban.soo.?aat
      បង់សូអ៊ិច
              ban.soo.?ic
      បង់សៀត
              ɓaŋ.siət
       បង់សែន ban.saen
បង់ស៊ីនប៉េនីស៊ីលីន
              ban.siin.pee.nii.sii.liin
       បង់ស៊ីស
              ban.siih
```







Phoneme Generator

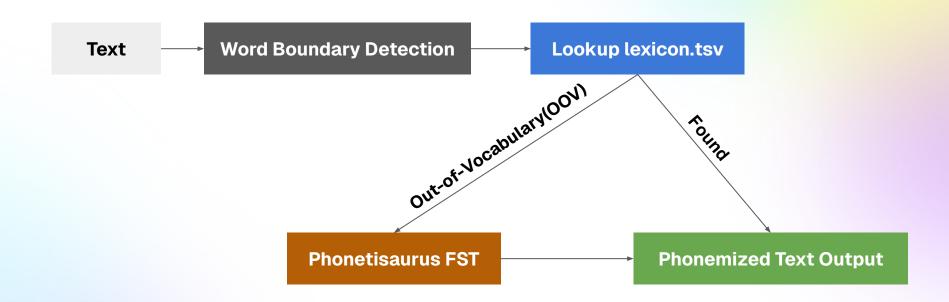








Text Phonemicization









Model Architecture Comparison

Architecture	Туре	Speed	Key Features	Strengths	Weaknesses		
Tacotron 2	Autoregressive seq2seq	Slow	Encoder-decoder with attention, generates mel-spectrograms	High naturalness, expressive speech, well-documented	Slow inference, can skip/repeat words, attention alignment issues		
FastSpeech / FastSpeech 2	Non-autoregressive feed-forward	Very Fast (270x faster than Tacotron)	Parallel generation, duration predictor, transformer-based	Extremely fast, robust (no word skipping), controllable speed	Requires teacher model or forced alignment for training		
VITS	End-to-end VAE-based	Fast	Variational autoencoder, monotonic alignment search, single-stage	Direct text-to-waveform, no separate vocoder needed, efficient	Complex architecture, GPU memory intensive		
Transformer TTS	Autoregressive transformer	Medium	Multi-head attention, transformer encoder-decoder	Better than RNN-based models, parallelizable training	Still slower than non-autoregressive models		







Dataset Requirements

- Clean Speech, No background noise, Consistent recording environment
- □ Sample Rate: 22.05 kHz \rightarrow 44.1 kHz
- Bit Depth: 16 bits per sample
- Channel: Mono
- Audio Format: WAV, FLAC (Lossless)
- Audio Max Duration: 10 seconds



Text អ្នកស្រី កែវ ចន្ទ បេះទុរេនលក់អស់៣០០តោនរួចហើយ ជាមួយតម្លៃ១៥០០០រៀលក្នុងមួយគីឡូក្រាម







Hardware Requirements

CPU: Good enough

GPU VRAM: 24GB Min.

RAM: 32GB

Disk: 1TB

Duration: A week or two (depends on dataset size)

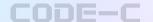






Knowledge Requirements

- Python / Virtualenv / Anaconda
- Machine Learning Concepts
- PyTorch
- Scipy
- Numpy
- C/C++ for binding
- etc.







Evaluation Metrics

MOS (Mean Opinion Score)

Mean Opinion Score (MOS) for TTS is a subjective evaluation method where human listeners rate synthesized speech on a scale (typically 1-5) for quality, naturalness, and clarity.

CMOS (Comparative Mean Opinion Score)

CMOS, or Comparative Mean Opinion Score, is a method for evaluating TTS quality by having human listeners directly compare two speech samples and choose which one is better.

Perceptual evaluation of speech quality (PESQ)

An objective, algorithmic method for assessing speech quality by comparing a degraded signal to a reference signal, making it applicable to Text-to-Speech (TTS) systems.







Deployment

TorchServe

ONNXRuntime by Microsoft (Anywhere)

NCNN from Tencent (Optimized for Mobile)

MNN from Alibaba (Optimized for Mobile)

Triton Inference Server by NVIDIA (NVIDIA Chips)





Demo

https://vok-tts-samples.netlify.app/

https://sarika.gov.kh







Next Challenges

- Khmer Word Tokenizer
- Neural Verbalizer
- Better Punctuation Model







References

https://github.com/seanghay/awesome-khmer-language

https://github.com/AdolfVonKleist/Phonetisaurus

https://github.com/seanghay/khmernormalizer

https://github.com/seanghay/khmercut

https://github.com/seanghay/tha

https://github.com/seanghay/phonetisaurus-js

https://github.com/sillsdev/khmer-character-specification

https://github.com/seanghay/automatic-phonemic-and-phonetic-transcription







Questions

Email: seanghay.dev@gmail.com

GitHub: @seanghay

Twitter: @seanghay_yath

Telegram:

https://t.me/seanghay_yath





