



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ

# Speech synthesis on a shoe string

Francis M. Tyers<sup>†</sup>

[ftyers@hse.ru](mailto:ftyers@hse.ru)

<https://www.hse.ru/org/persons/209454856>

Национальный исследовательский университет  
«Высшая школа экономики» (Москва)

†

with Josh Meyer, Jonathan North Washington and Marinyé Andreeva

**Introductions**

**Motivation**

**Project**

**Methodology**

# Introductions

*6 An endangered language will progress if its speakers can make use of electronic technology*

Main research interest:

- Language technology for maintenance and revival

I have worked with:

- Machine translation [Assimilation, Dissemination]
- Modelling morphology [Spellchecking]
- Syntactically-annotated corpora [Involve MLers]

And now:

- Speech synthesis

# **Motivation**

# Why speech synthesis?



- Interesting results with little data
- Lay groundwork for further speech-related work
- Real «application»

# Data requirements

## **Synthesis** (TTS)

- Hours of audio
- Studio quality
- Single speaker
- Not processor intensive

## **Recognition** (STT)


- 100s of hours of audio
- Telephone quality
- Many speakers
- Processor intensive


# Virtual assistants

- Talk to your computer
  - Information retrieval
  - «Smart» home stuff
- 3x faster than typing
- People really use them!
  - ...in English (or Russian)



# How they work





# **Project description**

# Project

## The idea:

- Collect a speech corpus of Chuvash
  - At least enough data to build a prototype
- Train a speech synthesis system
  - *Let's get the computer to speak Chuvash!*
- Develop an end-to-end “recipe”

## Budget:

- Recordings: 7000 RUB (approx. 100€)
  - 10€/hour
  - Estimated 1 min of audio = 2 mins recording
  - 10 minutes/day:
    - 5 hours audio (10h recording) in 1 month
- Train ticket (MSK → SpB): 3000 RUB (approx. 50€)

*Let's make a Chuvash voice!*  
Speech Synthesis Workshop, Higher School of Economics  
9th December, 2017




50€ well spent!

# Participants



Follow Josh on Twitter! @joshmeyerphd




- **Turkic language, Oghur sub-group**
- Spoken in Chuvashia, Volga region, Russia
- Approximately one million speakers (c. 2010)
- Intergenerational transmission breaking down

# Chuvash




- Turkic language, Oghur sub-group
- Spoken in Chuvashia, Volga region, Russia
- Approximately one million speakers (c. 2010)
- Intergenerational transmission breaking down

# Chuvash



- Turkic language, Oghur sub-group
- Spoken in Chuvashia, Volga region, Russia
- **Approximately one million speakers (c. 2010)**
- Intergenerational transmission breaking down

Распределение чувашей, говорящих на чувашском языке, в Чувашской Республике по возрасту\*



- Turkic language, Oghur sub-group
- Spoken in Chuvashia, Volga region, Russia
- Approximately one million speakers (c. 2010)
- Intergenerational transmission breaking down

«Most respondents, even Chuvash-speakers, think that Chuvash is not modern.»<sup>1</sup>

- Sociolinguistic study of pupils and teachers
- First in Chuvashia

---

<sup>1</sup>Alòs i Font, H. (2016) “The Chuvash language in the Chuvash Republic: An example of the rapid decline of one of Russia’s major languages”

# **Methodology**

# Collecting data

## **chuvas.org:**

- Online community
- Daily news in Chuvash
- Short articles 90-150 toks.
- CC-BY-SA licence

## **collection**

- wget
- BeautifulSoup



## Collecting data



[chuvash.org](http://chuvash.org)

- Online community
  - Daily news in Chuvash
  - Short articles 90-150 toks.
  - CC-BY-SA licence

## collection

- wget
  - BeautifulSoup



# Processing



Паша и Егор, сыновья бывшего каскадного самантра, Шупашкарти Мускав көпөрө синек инеке тухнә. Сына пұла транспорт сүресси сүр сөхеттөкө чарыннан пәттередүй. 23 сүлгі сәмбік водитель талғаса пыракан автобус сабд мани на ғына пырыса кінди.

Қайда Енри Шалти ғолен министерствин сайттөне леккін автобус 12-мөш маршрутта сүренинне айланып. Транспорта 5 сап ларса пынă, анчах вёсеменен инисін өтө аманманнинин ғыноттараңыз. Малтаның патайтын тұрақ, автобусын термөздө кимна пұттараң. Аның ку вересе хамында сирегелеттең иеттерек. Күзакансем қаланды тата жо дитеттің тәріхін тәріхін сиппа пәттөнметте майдан пур.

Инеке сәттінде спектаннистом палайтасынше ұнна тәржүмді. Сәй үнірексем тіріх каскадын, аларнан әмбәд машинасында халықтын күрнәлік.

Транспорт, ашыпашар, ғынносем, ашем, тайлірмаксем, күл-бір, автобусем, аМускав көпөрө

Хыныр сәнгек: <http://72.111.162.100/~1532922/>

1 0001 — Паян ирпе, қынсем ёче вакканың самантра, Шупашкарти Мускав көпөрө синек инеке тухнә. ↪

2 0002 — Җавна пұла транспорт сүресси сүр сөхеттөхе чарыннине пәлтересе. ↪

3 0003 — 23 сүлгі тәріхінде водитель тытса пыракан автобус сакарп машина үнне пырыса кінен. ↪

4 0004 — Чәваш Енри Шалти ёссен министерствин сайттөнде леккін автобус 12-мөш маршрутта сүренинде асанны. ↪

5 0005 — Транспорта 5 қын ларса пынă, анчах вёсеменен ніхаше те аманманнине ёнентере. ↪

- Bespoke sentence segmenter in Python
- Fixing minor encoding errors (Latin vs. Cyrillic)

# Audio recording materials



iPhone 4 S



- Recording device
- Noise insulation



Jonathan: « It's amazing that the hardest part is managing to get the cushions to stand up straight »

## **Recording:**

- Auphonic (for iPhone 4S)

## **Processing:**

- Google sheets
- sox - <http://sox.sourceforge.net/>
- pyAudioAnalysis -  
<https://github.com/tyiannak/pyAudioAnalysis>

## **Modelling:**

- Ossian -  
<https://github.com/CSTR-Edinburgh/Ossian>
- Merlin -  
<https://github.com/CSTR-Edinburgh/merlin>



# Recording

## Our instructions:

- 5 seconds of silence at the beginning
  - For the noise profile
- 2 seconds between each sentence
  - To facilitate splitting

## What Marinyé did:

- Produce a transcript
- Practice three times
- Read out numerals and abbreviations
- Correct errors in text — and note them down

Аса илтерер, Шулашқарти троллейбус управлениең үкәз-тәнің еңген ысырлайха көрсө үкім. Үйааш Еңбен Финанс министерстви предприятие Республика қызыннан 30 миллион тенге үйінә 0,1% ставкабы 3 сүлгілә пама паләрті. Ана үйк үйәнән 1-2-мәшесенне үйәрасшай.


4 Казахстан Президенти Нұрсултан Назарбаев паян Хүш кәләрді — пәнчәшерен те майыпен казах чөлхине латын саспаллисем қын күсарасси пирки.

Латын саспаллисем қын күсасси темиңе үтәйран тәрать. Чи маңтап — саспаллисем ыышне қырблептессы. Ку ең паянни хушула пурналдана та ёйті. Қене алғавитра пурб 32 саспалли, вәсемнен тәхжардамб — апострофлә пулә («ша», сәмахран, «с'» пек сыйынб).

Пәнчәшерле кириллицада 2025 үй тәліне хәтілма ынышыннан. Латын саспаллисем қын вәйләт тәліне көнеке пичетлессіне, пичет кәләрәмбесене, хут сарабайшын күсараса қитермелле. Шукулсанче латиницәла план тәрәх 2022 үлтандыра вәрентме тытәнб.

Латиница қын күсасси пирки Казахстанда 2012 үйлес калаңма тытәннан. Кәжалхы ака үйәкінчесе вара Н. Назарбаев латын саспаллисем қын күсмалы плана жатәрлеме ыйтты. Сәмәх май, Казахстанда латын саспаллийдесмене халық вәйләтре усқа курма пусанды та

# Data storage and organisation



- [uploadfiles.io](https://uploadfiles.io)
- Files last 30 days
- Easy to use


# Data storage and organisation

Recording project

ID	Текст	Слова	Запись	Дата	Длина	Качество	Комментарии	Шум
1	итого	28659			6:24:11			
2								
3								
4	<a href="https://chuvasch.org/news/17390.html">https://chuvasch.org/news/17390.html</a>	95	<a href="#">https://file.yolodou</a>	2017-11-10	0:01:05	✓		
5	<a href="https://chuvasch.org/news/17397.html">https://chuvasch.org/news/17397.html</a>	112	<a href="#">https://file.yod0m9</a>	2017-11-14	0:01:29			
6	<a href="https://chuvasch.org/news/17407.html">https://chuvasch.org/news/17407.html</a>	78	<a href="#">https://file.yolx3q7</a>	2017-11-14	0:01:21			
7	<a href="https://chuvasch.org/news/17394.html">https://chuvasch.org/news/17394.html</a>	108	<a href="#">https://file.yo939w9</a>	2017-11-14	0:01:34			
8	<a href="https://chuvasch.org/news/17362.html">https://chuvasch.org/news/17362.html</a>	92	<a href="#">https://file.yoxwds</a>	2017-11-14	0:01:20			
9	<a href="https://chuvasch.org/news/14981.html">https://chuvasch.org/news/14981.html</a>	97	<a href="#">https://file.yo9ggzg</a>	2017-11-14	0:01:17		high-amplitude white noise in the	
10	<a href="https://chuvasch.org/news/14951.html">https://chuvasch.org/news/14951.html</a>	170	<a href="#">https://file.yoylyq</a>	2017-11-15	0:01:55		high-amplitude white noise in the	
11	<a href="https://chuvasch.org/news/17389.html">https://chuvasch.org/news/17389.html</a>	91	<a href="#">https://file.yo9f7nd</a>	2017-11-14	0:01:20			
12	<a href="https://chuvasch.org/news/14459.html">https://chuvasch.org/news/14459.html</a>	94	<a href="#">https://file.yo9d5</a>	2017-11-15	0:01:14		high-amplitude white noise in the	
13	<a href="https://chuvasch.org/news/1793.html">https://chuvasch.org/news/1793.html</a>	97	<a href="#">https://file.yo9q9y</a>	2017-11-15	0:01:26			
14	<a href="https://chuvasch.org/news/16922.html">https://chuvasch.org/news/16922.html</a>	110	<a href="#">https://file.yo9ismay</a>	2017-11-18	0:01:28			
15	<a href="https://chuvasch.org/news/17289.html">https://chuvasch.org/news/17289.html</a>	89	<a href="#">https://file.yo9dn72</a>	2017-11-15	0:01:20			
16	<a href="https://chuvasch.org/news/17290.html">https://chuvasch.org/news/17290.html</a>	100	<a href="#">https://file.yo9um36</a>	2017-11-15	0:01:16			
17	<a href="https://chuvasch.org/news/17414.html">https://chuvasch.org/news/17414.html</a>	106	<a href="#">https://file.yo9of48</a>	2017-11-15	0:01:16			
18	<a href="https://chuvasch.org/news/17409.html">https://chuvasch.org/news/17409.html</a>	114	<a href="#">https://file.yo9phth</a>	2017-11-15	0:01:37			


- Links files with text
- Some metadata
- Comments about quality, edits

# Processing



- Use `sox` to make a silence profile
- Clean the audio using the profile
- Trim the first three seconds

# Processing



- Detect sentence spans using pyAudioAnalysis
- Split sentences using sox.

Result:



Çемёрлесем хула тăрăх пĕр каяччă  
хайне евĕр транспортпа çўренине  
асăрханă.



Вăл гироскутера аса илтернĕ.



Иртен-çўрен ёна видо ўкерсе халăх те-  
телёсене вырнаçтарнă.

- Text and sound sentence aligned corpus

	<b>Files</b>	<b>Sentences</b>	<b>Words</b>	<b>Minutes</b>
<b>Total:</b>	299	2994	34648	384 (6h24m)
<b>Postproc:</b>	226	2090	21815	200 (3h20m)

### **Loss:**

- 5 seconds per file (approx. 24m)
- 2 seconds per sentence (approx. 1h40m)  
= 2h04m.
- Max: 4h20m
- minus approx. 25% of files
- → 3h20m

# Modelling


Шупашкартан Питёре  
вёсес килет-и?

«Do you want to fly from  
Şupuşkar to St. Petersburg?»



RAW TEXT

Шупашкортан Питёре  
вёсес килет-и?



# Modelling




Шупашкортан Питёре  
вёсес килет-и?



# Modelling

Шупашкортан Питёре  
вёсес килет-и?




# Modelling

Шупашкортан Питёре  
вёсес килет-и?



**HARD**






**HARD**

# Modelling


Шупашкортан Питёре  
вёсес килет-и?




**EASIER!**

Шупашкортан Питёре  
вёсес килет-и?

**REGRESSION**



# Modelling



# Modelling

Шупашкортан Питёре  
вёсес килет-и?



[00100000070002301000000000900]

[6010000000000008]



# FRONTEND

## *Ossian*

Ossian




Шупашкортан Питĕре  
вĕçес килет-и?

[001000007000230100000000900]

# BACKEND

## *Merlin*

# *Merlin*



[6010000000000008]



# Frontend: Text features

*Шупашкартан Питёре вёçес килет-и?*

## Letters

ш у п а ш к а р т а н п и т ё р е в ё ç е с к и л е т - и ?

## Phonemes

ʂ u b a ş k a r d a n p i t ə r ε v ə z ε s k i l ε t \_ i \_

## Words

[ʂ u b a ş k a r d a n] [p i t ə r ε] [v ə z ε s] [k i l ε t] [\_ i ]

## Parts of speech


[PROPN ʂ u b a ş k a r d a n] [PROPN p i t ə r ε] [VERB v ə z ε s] [VERB k i l ε t] [AUX \_ i \_]

## Syllables


[PROPN [ʂ u] [b a ş] [k a r] [d a n]] [PROPN [p i] [t ə] [r ε]] [VERB [v ə z] [ε s]] [VERB [k i l] [ε t]] [AUX [\_ i \_]]

etc.


# Frontend: Features to segments



# Frontend: Features to segments



# Frontend: Features to segments




c-1:n, c0:p, c+1:i, w0:Питёре, w-1:Шупашкартан, w+1:вёces, s-1:dan, s+1:t□, ...



[ 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 1 0 ]

# Backend: Audio features



- Each feature vector is a frame of audio (10ms)
- Features are amplitudes at frequency ranges

## Text feature vectors

```
[0000010000000000001000080000500100000000000]  
[0000010040000000001000070000000100000000000]  
[0000010000500000001000000009000100000000000]
```




## Audio feature vectors

- Unequal number of vectors for text and audio
  - ...so just repeat them until they align

# Training


```
[00000100000000000010000800005001000000000000  
[00000100000000000010000800005001000000000000  
[00000100000000000010000800005001000000000000  
[00000100000000000010000800005001000000000000  
[00000100000000000010000800005001000000000000]
```



```
[0010000100000000  
[002000300070002  
[301000000000900  
[005090010004000  
[6010000100000008]
```

# Training


```
[0000001000000000000000001000008000050010000000000000000  
[0000001000000000000000001000008000050010000000000000000  
[0000001000000000000000001000008000050010000000000000000  
[0000001000000000000000001000008000050010000000000000000  
[0000001000000000000000001000008000050010000000000000000]
```



```
[0010000100000000  
[002000300070002  
[3010000000009001  
[0050900100040000  
[6010000100000008]
```

# Training


```
[0000001000000000000010000800005001000000000000]  
[0000001000000000000010000800005001000000000000]  
[0000001000000000000010000800005001000000000000]  
[0000001000000000000010000800005001000000000000]
```



```
[0010000100000000]  
[002000300070002]  
[3010000000009001]  
[0050900100040001]  
[6010000100000008]
```

# Training

```
[0000001000000000000010000800005001000000000000]  
[0000001000000000000010000800005001000000000000]  
[0000001000000000000010000800005001000000000000]  
[0000001000000000000010000800005001000000000000]  
[0000001000000000000010000800005001000000000000]
```



```
[0010000100000000]  
[002000300070002]  
[301000000000900]  
[005090010004000]  
[601000010000008]
```

# Evaluation

- 2 evaluators
- Scores 1 (very bad) — 5 (very good)

	Russian*	Chuvash
<b>Intelligibility</b>	5, 5	5, 3
<b>Pronunciation</b>	5, 5	5, 3
<b>Stress</b>	5, 5	4, 3
<b>Intonation</b>	2, 5	2, 3

- Intelligible
- Pronunciation and stress mostly ok
- Flat intonation

\* Google TTS

# **Discussion**

# Future directions


## **Speech synthesis**

- Better linguistic features [POS, morphs, etc.]
- Fix segmentation of remaining 25%
- Numeral and abbreviation processing

## **Speech recognition**

- Systems
  - Kaldi
  - Mozilla DeepSpeech
- Collecting more data
  - More volume (100s of hours)
    - How ?
  - More varied (different recording conditions)

# Common Voice



Voice is natural, voice is human. That's why we're fascinated with creating usable voice technology for our machines. But to create voice systems, an extremely large amount of voice data is required.

Most of the data used by large companies isn't available to the majority of people. We think that stifles innovation. So we've launched Project Common Voice, a project to help make voice recognition open to everyone.

[READ MORE](#)

- Web platform for collecting voice data
- Requires CC-0 — mostly impossible
- But...
  - Backend software is free/open-source



**Ман қине вайхат уйарнишён тавах!**

*Merci beaucoup pour votre attention!*

[http://www.github.com/ftyers/Turkic\\_TTS](http://www.github.com/ftyers/Turkic_TTS)